# Error analysis of the s-step Lanczos method in finite precision

*Erin Carson*
*James Demmel*

Electrical Engineering and Computer Sciences
University of California at Berkeley

May 6, 2014

| | | Form Approved OMB No. 0704-0188 |
|---|---|---|

**Report Documentation Page**

| 1. REPORT DATE **06 MAY 2014** | 2. REPORT TYPE | 3. DATES COVERED **00-00-2014 to 00-00-2014** |
|---|---|---|

| 4. TITLE AND SUBTITLE **Error analysis of the s-step Lanczos method in finite precision** | 5a. CONTRACT NUMBER |
|---|---|
| | 5b. GRANT NUMBER |
| | 5c. PROGRAM ELEMENT NUMBER |
| 6. AUTHOR(S) | 5d. PROJECT NUMBER |
| | 5e. TASK NUMBER |
| | 5f. WORK UNIT NUMBER |
| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) **University of California at Berkeley,Electrical Engineering and Computer Sciences,Berkeley,CA,94720** | 8. PERFORMING ORGANIZATION REPORT NUMBER |
| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | 10. SPONSOR/MONITOR'S ACRONYM(S) |
| | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) |

12. DISTRIBUTION/AVAILABILITY STATEMENT
**Approved for public release; distribution unlimited**

13. SUPPLEMENTARY NOTES

14. ABSTRACT
**The s-step Lanczos method is an attractive alternative to the classical Lanczos method as it enables an O(s) reduction in data movement over a xed number of iterations. This can signi cantly improve performance on modern computers. In order for s-step methods to be widely adopted, it is important to better understand their error properties. Although the s-step Lanczos method is equivalent to the classical Lanczos method in exact arithmetic, empirical observations demonstrate that it can behave quite di erently in nite precision. In the s-step Lanczos method the computed Lanczos vectors can lose orthogonality at a much quicker rate than the classical method a property which seems to worsen with increasing s. In this paper, we present, for the rst time, a complete rounding error analysis of the s-step Lanczos method. Our methodology is analogous to Paige's rounding error analysis for the classical Lanczos method [IMA J. Appl. Math., 18(3):341{349, 1976]. Our analysis gives upper bounds on the loss of normality of and orthogonality between the computed Lanczos vectors, as well as a recurrence for the loss of orthogonality. The derived bounds are very similar to those of Paige for classical Lanczos, but with the addition of an ampli cation term which depends on the condition number of the Krylov bases computed every s-steps. Our results con rm theoretically what is well-known empirically: the conditioning of the Krylov bases plays a large role in determining nite precision behavior.**

15. SUBJECT TERMS

| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT **unclassified** | b. ABSTRACT **unclassified** | c. THIS PAGE **unclassified** | **Same as Report (SAR)** | **29** | |

# ERROR ANALYSIS OF THE $S$-STEP LANCZOS METHOD IN FINITE PRECISION

ERIN CARSON AND JAMES DEMMEL

**Abstract.** The $s$-step Lanczos method is an attractive alternative to the classical Lanczos method as it enables an $O(s)$ reduction in data movement over a fixed number of iterations. This can significantly improve performance on modern computers. In order for $s$-step methods to be widely adopted, it is important to better understand their error properties. Although the $s$-step Lanczos method is equivalent to the classical Lanczos method in exact arithmetic, empirical observations demonstrate that it can behave quite differently in finite precision. In the $s$-step Lanczos method, the computed Lanczos vectors can lose orthogonality at a much quicker rate than the classical method, a property which seems to worsen with increasing $s$.

In this paper, we present, for the first time, a complete rounding error analysis of the $s$-step Lanczos method. Our methodology is analogous to Paige's rounding error analysis for the classical Lanczos method [*IMA J. Appl. Math.*, 18(3):341–349, 1976]. Our analysis gives upper bounds on the loss of normality of and orthogonality between the computed Lanczos vectors, as well as a recurrence for the loss of orthogonality. The derived bounds are very similar to those of Paige for classical Lanczos, but with the addition of an amplification term which depends on the condition number of the Krylov bases computed every $s$-steps. Our results confirm theoretically what is well-known empirically: the conditioning of the Krylov bases plays a large role in determining finite precision behavior.

**Key words.** Krylov subspace methods, error analysis, finite precision, roundoff, Lanczos, avoiding communication, orthogonal bases

**AMS subject classifications.** 65G50, 65F10, 65F15, 65N15, 65N12

**1. Introduction.** Given an $n$-by-$n$ symmetric matrix $A$ and a starting vector $v_0$ with unit 2-norm, $m$ steps of the Lanczos method [21] theoretically produces the orthonormal matrix $V_m = [v_0, \dots, v_m]$ and the $(m+1)$-by-$(m+1)$ symmetric tridiagonal matrix $T_m$ such that

$$AV_m = V_m T_m + \beta_{m+1} v_{m+1} e_{m+1}^T. \tag{1.1}$$

When $m = n - 1$, the eigenvalues $T_{n-1}$ are the eigenvalues of $A$. In practice, the eigenvalues of $T$ are still good approximations to the eigenvalues of $A$ when $m \ll n-1$, which makes the Lanczos method attractive as an iterative procedure. Many Krylov subspace methods (KSMs), including those for solving linear systems and least squares problems, are based on the Lanczos method. In turn, these various Lanczos-based methods are the core components in numerous scientific applications.

Classical implementations of Krylov methods, the Lanczos method included, require one or more sparse matrix-vector multiplications (SpMVs) and one or more inner product operations in each iteration. These computational kernels are both communication-bound on modern computer architectures. To perform an SpMV, each processor must communicate entries of the source vector it owns to other processors in the parallel algorithm, and in the sequential algorithm the matrix $A$ must be read from slow memory. Inner products involve a global reduction in the parallel algorithm, and a number of reads and writes to slow memory in the sequential algorithm (depending on the size of the vectors and the size of the fast memory).

Thus, many efforts have focused on communication-avoiding Krylov subspace methods (CA-KSMs), or $s$-step Krylov methods, which can perform $s$ iterations with $O(s)$ less communication than classical KSMs; see, e.g., [4, 5, 7, 9, 10, 16, 17, 8, 33, 35]. In practice, this can translate into significant speedups for many problems [24].

Equally important to the performance of each iteration is the convergence rate of the method, i.e., the total number of iterations required until the desired convergence criterion is met. Although theoretically the Lanczos process described by (1.1) produces an orthogonal basis and a tridiagonal matrix similar to $A$ after $n$ steps, these properties need not hold in finite precision. The effects of roundoff error on the ideal Lanczos process were known to Lanczos when he published his algorithm in 1950. Since then, much research has been devoted to better understanding this behavior, and to devise more robust and stable algorithms.

Although $s$-step Krylov methods are mathematically equivalent to their classical counterparts in exact arithmetic, it perhaps comes as no surprise that their finite precision behavior may differ significantly, and that the theories developed for classical methods in finite precision do not hold for the $s$-step case. It has been empirically observed that the behavior of $s$-step Krylov methods deviates further from that of the classical method as $s$ increases, and that the severity of this deviation is heavily influenced by the polynomials used for the $s$-step Krylov bases (see, e.g., [1, 4, 17, 18]).

Arguably the most revolutionary work in the finite precision analysis of classical Lanczos was a series of papers published by Paige [25, 26, 27, 28]. Paige's analysis succinctly describes how rounding errors propagate through the algorithm to impede orthogonality. These results were developed to give theorems which link the loss of orthogonality to convergence of the computed eigenvalues [28]. No analogous theory currently exists for the $s$-step Lanczos method.

In this paper, we present, for the first time, a complete rounding error analysis of the $s$-step Lanczos method. Our analysis here for $s$-step Lanczos closely follows Paige's rounding error analysis for orthogonality in classical Lanczos [27].

We present upper bounds on the normality of and orthogonality between the computed Lanczos vectors, as well as a recurrence for the loss of orthogonality. The derived bounds are very similar to those of Paige for classical Lanczos, but with the addition of an amplification term which depends on the condition number of the Krylov bases computed every $s$ steps. Our results confirm theoretically what is well-known empirically: the conditioning of the Krylov bases plays a large role in determining finite precision behavior. In particular, if one can guarantee that the basis condition number is not too large throughout the iteration, the loss of orthogonality in the $s$-step Lanczos method should not be too much worse than in classical Lanczos. As Paige's subsequent groundbreaking convergence analysis [28] was based largely on the results in [27], our analysis here similarly serves as a stepping stone to further understanding of the $s$-step Lanczos method.

The remainder of this paper is outlined as follows. In Section 2, we present related work in $s$-step Krylov methods and the analysis of finite precision Lanczos. In Section 3, we review a variant of the Lanczos method and derive the corresponding $s$-step Lanczos method, as well as provide a numerical example that will help motivate our analysis. In Section 4, we first state our main result in Theorem 4.2 and comment on its interpretation; the rest of the section is devoted to its proof. In Section 5, we recall our numerical example from Section 3 in order to demonstrate the bounds proved in Section 4. Section 6 concludes with a discussion of future work.

**2. Related work.** We briefly review related work in $s$-step Krylov methods as well as work related to the analysis of classical Krylov methods in finite precision.

**2.1. $s$-step Krylov subspace methods.** The term '$s$-step Krylov method', first used by Chronopoulos and Gear [6], describes variants of Krylov methods where the iteration loop is split into blocks of $s$ iterations. Since the Krylov subspaces

required to perform $s$ iterations of updates are known, bases for these subspaces can be computed upfront, inner products between basis vectors can be computed with one block inner product, and then $s$ iterations are performed by updating the coordinates in the generated Krylov bases (see Section 3 for details). Many formulations and variations have been derived over the past few decades with various motivations, namely increasing parallelism (e.g., [6, 35, 36]) and avoiding data movement, both between levels of the memory hierarchy in sequential methods and between processors in parallel methods. A thorough treatment of related work can be found in [17].

Many empirical studies of $s$-step Krylov methods found that convergence often deteriorated using $s > 5$ due to the inherent instability of the monomial basis. This motivated research into the use of better-conditioned bases (e.g., Newton or Chebyshev polynomials) for the Krylov subspace, which allowed convergence for higher $s$ values (see, e.g., [1, 16, 18, 31]). Hoemmen has used a novel matrix equilibration and balancing approach to achieve similar effects [17].

The term 'communication-avoiding Krylov methods' refers to $s$-step Krylov methods and implementations which aim to improve performance by asymptotically decreasing communication costs, possibly both in computing inner products and computing the $s$-step bases, for both sequential and parallel algorithms; see [9, 17]. Hoemmen et al. [17, 24] derived communication-avoiding variants of Lanczos, Arnoldi, Conjugate Gradient (CG) and the Generalized Minimum Residual Method (GMRES). Details of nonsymmetric Lanczos-based CA-KSMs, including communication-avoiding versions of Biconjugate Gradient (BICG) and Stabilized Biconjugate Gradient (BICGSTAB) can be found in [4]. Although potential performance improvement is our primary motivation for studying these methods, we use the general term '$s$-step methods' here as our error analysis is independent of performance.

Many efforts have been devoted specifically to the $s$-step Lanczos method. The first $s$-step Lanczos methods known in the literature are due to Kim and Chronopoulos, who derived a three-term symmetric $s$-step Lanczos method [19] as well as a three-term nonsymmetric $s$-step Lanczos method [20]. Hoemmen derived a three-term communication-avoiding Lanczos method, CA-Lanczos [17]. Although the three-term variants require less memory, their numerical accuracy can be worse than implementations which use two coupled two-term recurrences [15]. A two-term communication-avoiding nonsymmetric Lanczos method (called CA-BIOC, based on the 'BIOC' version of nonsymmetric Lanczos of Gutknecht [14]) can be found in [2]. This work includes the derivation of a new version of the $s$-step Lanczos method, equivalent in exact arithmetic to the variant used by Paige [27]. It uses a two-term recurrence like BIOC, but is restricted to the symmetric case and uses a different starting vector.

For $s$-step KSMs that solve linear systems, increased roundoff error in finite precision can decrease the maximum attainable accuracy of the solution, resulting in a less accurate solution than found by the classical method. A quantitative analysis of roundoff error in CA-CG and CA-BICG can be found in [3]. Based on the work of [34] for conventional KSMs, we have also explored implicit residual replacement strategies for CA-CG and CA-BICG as a method to limit the deviation of true and computed residuals when high accuracy is required (see [3]).

**2.2. Error analysis of the Lanczos method.** Lanczos and others recognized early on that rounding errors could cause the Lanczos method to deviate from its ideal theoretical behavior. Since then, various efforts have been devoted to analyzing, and explaining, and improving the finite precision Lanczos method.

Widely considered to be the most significant development was the series of pa-

pers by Paige discussed in Section 1. Another important development was due to
Greenbaum and Strakoš, who performed a backward-like error analysis which showed
that finite precision Lanczos and CG behave very similarly to the exact algorithms
applied to any of a certain class of larger matrices [12]. Paige has recently shown
a similar type of augmented stability for the Lanczos process [29]. There are many
other analyses of the behavior of various KSMs in finite precision, including some
more recent results due to Wülling [37] and Zemke [38]; for a thorough overview of
the literature, see [22, 23].

A number of strategies for maintaining the orthogonality among the Lanczos
vectors were inspired by the analysis of Paige, such as selective reorthogonalization [30]
and partial reorthogonalization [32]. Recently, Gustafsson et al. have extended such
reorthogonalization strategies for classical Lanczos to the $s$-step case [13].

**3. The $s$-step Lanczos method.** The classical Lanczos method is shown in
Algorithm 1. We use the same variant of Lanczos as used by Paige in his error
analysis for classical Lanczos [27] to allow easy comparison of results. This is the
first instance of an $s$-step version of this particular Lanczos variant; other existing
$s$-step Lanczos variants are described in Section 2.1. Note that as in [27] our analysis
will assume no breakdown occurs and thus breakdown conditions are not discussed
here. We now give a derivation of $s$-step Lanczos, obtained from classical Lanczos in
Algorithm 1.

---

**Algorithm 1** Lanczos

---

**Require:** $n$-by-$n$ real symmetric matrix $A$ and length-$n$ starting vector $v_0$ such that
 $\|v_0\|_2 = 1$
 1: $u_0 = Av_0$
 2: **for** $m = 0, 1, \ldots$ until convergence **do**
 3:     $\alpha_m = v_m^T u_m$
 4:     $w_m = u_m - \alpha_m v_m$
 5:     $\beta_{m+1} = \|w_m\|_2$
 6:     $v_{m+1} = w_m / \beta_{m+1}$
 7:     $u_{m+1} = Av_{m+1} - \beta_{m+1} v_m$
 8: **end for**

---

Suppose we are beginning iteration $m = sk$ where $k \in \mathbb{N}$ and $0 < s \in \mathbb{N}$. By
induction on lines 6 and 7 of Algorithm 1, we can write

$$v_{sk+j}, u_{sk+j} \in \mathcal{K}_{s+1}(A, v_{sk}) + \mathcal{K}_{s+1}(A, u_{sk}) \qquad (3.1)$$

for $j \in \{0, \ldots, s\}$, where $\mathcal{K}_i(A, x) = \mathrm{span}\{x, Ax, \ldots, A^{i-1}x\}$ denotes the Krylov sub-
space of dimension $i$ of matrix $A$ with respect to vector $x$. Note that since $u_0 = Av_0$,
if $k = 0$ we have

$$v_j, u_j \in \mathcal{K}_{s+2}(A, v_0).$$

for $j \in \{0, \ldots, s\}$.

For $k > 0$, we then define 'basis matrix' $\mathcal{Y}_k = [\mathcal{V}_k, \mathcal{U}_k]$, where $\mathcal{V}_k$ and $\mathcal{U}_k$ are size
$n$-by-$(s+1)$ matrices whose columns form bases for $\mathcal{K}_{s+1}(A, v_{sk})$ and $\mathcal{K}_{s+1}(A, u_{sk})$,
respectively. For $k = 0$, we define $\mathcal{Y}_0$ to be a size $n$-by-$(s+2)$ matrix whose columns
form a basis for $\mathcal{K}_{s+2}(A, v_0)$. Then by (3.1), we can represent $v_{sk+j}$ and $u_{sk+j}$, for
$j \in \{0, \ldots, s\}$, by their coordinates (denoted with primes) in $\mathcal{Y}_k$, i.e.,

$$v_{sk+j} = \mathcal{Y}_k v'_{k,j}, \qquad u_{sk+j} = \mathcal{Y}_k u'_{k,j}. \qquad (3.2)$$

Note that for $k = 0$, the coordinate vectors are length $s + 2$ and for $k > 0$, the coordinate vectors are length $2s + 2$. We can write a similar equation for auxiliary vector $w_{sk+j}$, i.e., $w_{sk+j} = \mathcal{Y}_k w'_{k,j}$ for $j \in \{0, \dots, s-1\}$. We define also the Gram matrix $G_k = \mathcal{Y}_k^T \mathcal{Y}_k$, which is size $(s+2)$-by-$(s+2)$ for $k = 0$ and $(2s+2)$-by-$(2s+2)$ for $k > 0$. Using this matrix, the inner products in lines 3 and 5 can be written

$$\alpha_{sk+j} = v_{sk+j}^T u_{sk+j} = v_{k,j}'^T \mathcal{Y}_k^T \mathcal{Y}_k u'_{k,j} = v_{k,j}'^T G_k u'_{k,j} \quad \text{and} \tag{3.3}$$

$$\beta_{sk+j+1} = (w_{sk+j}^T w_{sk+j})^{1/2} = (w_{k,j}'^T \mathcal{Y}_k^T \mathcal{Y}_k w'_{k,j})^{1/2} = (w_{k,j}'^T G_k w'_{k,j})^{1/2}. \tag{3.4}$$

We assume that the bases are generated via polynomial recurrences represented by the matrix $\mathcal{B}_k$, which is in general upper Hessenberg but often tridiagonal in practice. The recurrence can thus be written in matrix form as

$$A\hat{\underline{\mathcal{Y}}}_k = \hat{\mathcal{Y}}_k \mathcal{B}_k$$

where $\mathcal{B}_k$ is size $(s+2)$-by-$(s+2)$ for $k = 0$ and size $(2s+2)$-by-$(2s+2)$ for $k > 0$, and $\hat{\underline{\mathcal{Y}}}_k = \left[ \hat{\mathcal{V}}_k[I_s, 0_{s,1}]^T, 0_{n,1}, \hat{\mathcal{U}}_k[I_s, 0_{s,1}]^T, 0_{n,1} \right]$. Therefore, for $j \in \{0, \dots, s-1\}$,

$$Av_{sk+j+1} = A\mathcal{Y}_k v'_{k,j+1} = A\hat{\underline{\mathcal{Y}}}_k v'_{k,j+1} = \mathcal{Y}_k \mathcal{B}_k v'_{k,j+1}. \tag{3.5}$$

Thus, to compute iterations $sk + 1$ through $sk + s$ in $s$-step Lanczos, we first generate basis matrix $\mathcal{Y}_k$ such that (3.5) holds, and we compute the Gram matrix $G_k$ from the resulting basis matrix. Then updates to the length-$n$ vectors can be performed by updating instead the length-$(2s + 2)$ coordinates for those vectors in $\mathcal{Y}_k$. Inner products and multiplications with $A$ become smaller operations which can be performed locally, as in (3.3), (3.4), and (3.5). The complete $s$-step Lanczos algorithm is presented in Algorithm 2. Note that in Algorithm 2 we show the length-$n$ vector updates in each inner iteration (lines 16 and 18) for clarity, although these vectors play no part in the inner loop iteration updates. In practice, the basis change operation (3.2) can be performed on a block of coordinate vectors at the end of each outer loop to recover $v_{sk+i}$ and $u_{sk+i}$, for $i \in \{1, \dots, s\}$.

**3.1. A numerical example.** We give a brief example to demonstrate the behavior of $s$-step Lanczos in finite precision and to motivate our theoretical analysis. We run $s$-step Lanczos (Algorithm 2) on a 2D Poisson matrix with $n = 256$, $\|A\|_2 = 7.93$, using a random starting vector. The same starting vector is used in all tests, which were run using double precision. Results for classical Lanczos run on the same problem are shown in Figure 3.1 for comparison. In Figure 3.2, we show $s$-step Lanczos results for $s = 2$ (left), $s = 4$ (middle), and $s = 8$ (right), using monomial (top), Newton (middle), and Chebyshev (top) polynomials for computing the bases in line 3. The plots show the number of eigenvalue estimates (Ritz values) that have converged, within some relative tolerance, to a true eigenvalue over the iterations. Note that we do not count duplicates, i.e., multiple Ritz values that have converged to the same eigenvalue of $A$. The solid black line $y = x$ represents the upper bound.

From Figure 3.2 we see that for $s = 2$, $s$-step Lanczos with the monomial, Newton, and Chebyshev bases all well-replicate the convergence behavior of classical Lanczos; for the Chebyshev basis the plots look almost identical. However, as $s$ increases, we see that both convergence rate and accuracy to which we can find approximate eigenvalues within $n$ iterations decreases for all bases. This is clearly the most drastic for the monomial basis case; e.g., for the Chebyshev and Newton bases with $s = 8$,

---

**Algorithm 2** $s$-step Lanczos

---

**Require:** $n$-by-$n$ real symmetric matrix $A$ and length-$n$ starting vector $v_0$ such that
$\|v_0\|_2 = 1$

 1: $u_0 = Av_0$
 2: **for** $k = 0, 1, \ldots$ until convergence **do**
 3:     Compute $\mathcal{Y}_k$ with change of basis matrix $\mathcal{B}_k$
 4:     Compute $G_k = \mathcal{Y}_k^T \mathcal{Y}_k$
 5:     $v'_{k,0} = e_1$
 6:     **if** $k = 0$ **then**
 7:         $u'_{0,0} = \mathcal{B}_k e_1$
 8:     **else**
 9:         $u'_{k,0} = e_{s+2}$
10:     **end if**
11:     **for** $j = 0, 1, \ldots, s-1$ **do**
12:         $\alpha_{sk+j} = v'^T_{k,j} G_k u'_{k,j}$
13:         $w'_{k,j} = u'_{k,j} - \alpha_{sk+j} v'_{k,j}$
14:         $\beta_{sk+j+1} = (w'^T_{k,j} G_k w'_{k,j})^{1/2}$
15:         $v'_{k,j+1} = w'_{k,j}/\beta_{sk+j+1}$
16:         $v_{sk+j+1} = \mathcal{Y}_k v'_{k,j+1}$
17:         $u'_{k,j+1} = \mathcal{B}_k v'_{k,j+1} - \beta_{sk+j+1} v'_{k,j}$
18:         $u_{sk+j+1} = \mathcal{Y}_k u'_{k,j+1}$
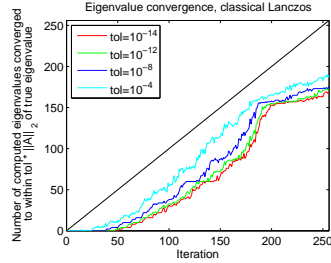19:     **end for**
20: **end for**

---



Fig. 3.1. *Number of converged Ritz values versus iteration number for classical Lanczos.*

we can at least still find eigenvalues to within relative accuracy $\sqrt{\epsilon}$ at the same rate as the classical case.

It is clear that the choice of basis used to generate Krylov subspaces affects the behavior of the method in finite precision. Although this is well-studied empirically in the literature, many theoretical questions remain open about exactly how, where, and to what extent the properties of the bases affect the method's behavior. Our analysis is a significant step toward addressing these questions.

**4. The $s$-step Lanczos method in finite precision.** Throughout our analysis, we use a standard model of floating point arithmetic where we assume the computations are carried out on a machine with relative precision $\epsilon$ (see [11]). Throughout the analysis we ignore $\epsilon$ terms of order $> 1$, which have negligible effect on our results. We also ignore underflow and overflow. Following Paige [27], we use the $\epsilon$ symbol to represent the relative precision as well as terms whose absolute values are bounded
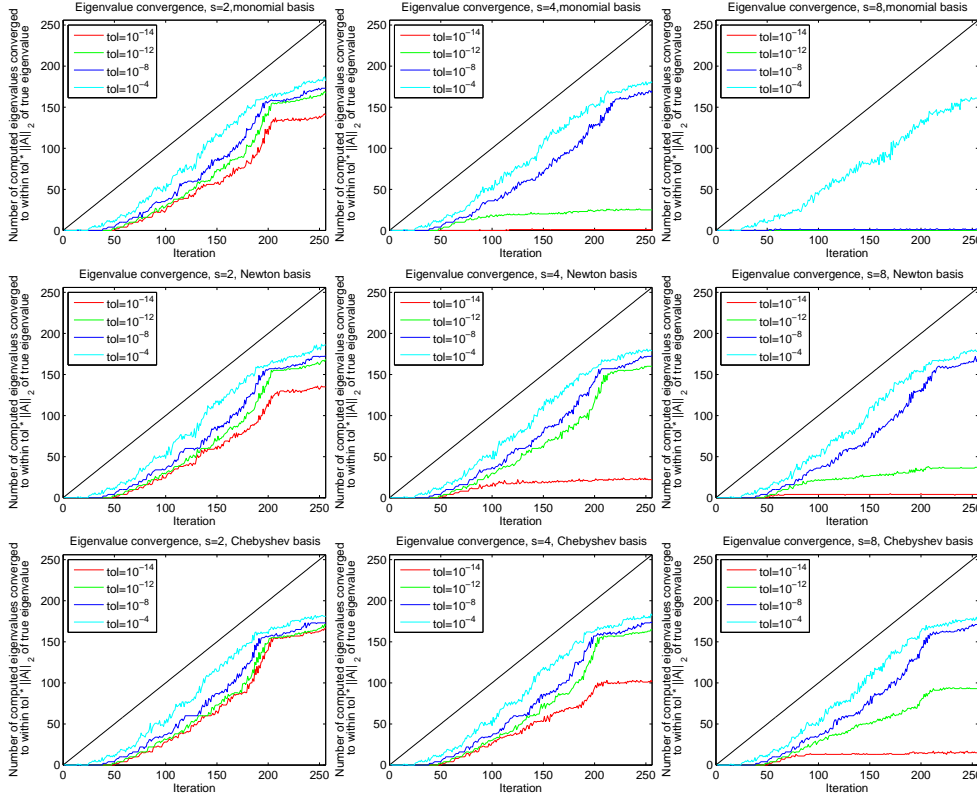
FIG. 3.2. *Number of converged Ritz values versus iteration number for s-step Lanczos using monomial (top), Newton (middle), and Chebyshev (bottom) bases for $s = 2$ (left), $s = 4$ (middle), and $s = 8$ (right).*

by the relative precision.

We will model floating point computation using the following standard conventions (see, e.g., [11]): for vectors $u, v \in \mathbb{R}^n$, matrices $A \in \mathbb{R}^{n \times m}$ and $G \in \mathbb{R}^{n \times n}$, and scalar $\alpha$,

$$
\begin{aligned}
fl(u - \alpha v) &= u - \alpha v - \delta w, & |\delta w| &\leq (|u| + 2|\alpha v|)\epsilon, \\
fl(v^T u) &= (v + \delta v)^T u, & |\delta v| &\leq n\epsilon |v|, \\
fl(Au) &= (A + \delta A)u, & |\delta A| &\leq m\epsilon |A|, \\
fl(A^T A) &= A^T A + \delta E, & |\delta E| &\leq n\epsilon |A^T||A|, \quad \text{and} \\
fl(u^T(Gv)) &= (u + \delta u)^T(G + \delta G)v, & |\delta u| &\leq n\epsilon |u|, |\delta G| \leq n\epsilon |G|.
\end{aligned}
$$

where $fl()$ represents the evaluation of the given expression in floating point arithmetic and terms with $\delta$ denote error terms. We decorate quantities computed in finite precision arithmetic with hats, e.g., if we are to compute the expression $\alpha = v^T u$ in finite precision, we get $\hat{\alpha} = fl(v^T u)$.

We first prove the following lemma, which will be useful in our analysis.

LEMMA 4.1. *Assume we have rank-r matrix $Y \in \mathbb{R}^{n \times r}$, where $n \geq r$. Let $Y^+$ denote the pseudoinverse of $Y$, i.e., $Y^+ = (Y^T Y)^{-1} Y^T$. Then for any vector $x \in \mathbb{R}^r$,*

*we can bound*

$$\| |Y| |x| \|_2 \leq \| |Y| \|_2 \|x\|_2 \leq \Gamma \|Yx\|_2.$$

*where* $\Gamma = \left\|Y^+\right\|_2 \| |Y| \|_2 \leq \sqrt{r} \left\|Y^+\right\|_2 \|Y\|_2$.
    *Proof.* We have

$$\| |Y||x| \|_2 \leq \| |Y| \|_2 \|x\|_2 \leq \| |Y| \|_2 \|Y^+ Y x\|_2 \leq \| |Y| \|_2 \|Y^+\|_2 \|Yx\|_2 \leq \Gamma \|Yx\|_2.$$

▢

    We note that the term $\Gamma$ can be thought of as a type of condition number for the matrix $Y$. In the analysis, we will apply the above lemma to the computed 'basis matrix' $\hat{\mathcal{Y}}_k$. We assume throughout that the generated bases $\hat{\mathcal{U}}_k$ and $\hat{\mathcal{V}}_k$ are numerically full rank. That is, all singular values of $\hat{\mathcal{U}}_k$ and $\hat{\mathcal{V}}_k$ are greater than $\epsilon n \cdot 2^{\lfloor \log_2 \sigma_1 \rfloor}$ where $\sigma_1$ is the largest singular value of $A$. The results of this section are summarized in the following theorem:

    THEOREM 4.2. *Assume that Algorithm 2 is implemented in floating point with relative precision $\epsilon$ and applied for $sk+j$ steps to the $n$-by-$n$ real symmetric matrix $A$, starting with vector $v_0$ with $\|v_0\|_2 = 1$. Let $\sigma = \||A|\|_2 / \|A\|_2$ and $\tau_k = \||\mathcal{B}_k|\|_2 / \|A\|_2$, where $\mathcal{B}_k$ is defined in (3.5), and let*

$$\bar{\Gamma}_k = \max_{i \in \{0,\ldots,k\}} \|\hat{\mathcal{Y}}_i^+\|_2 \| |\hat{\mathcal{Y}}_i| \|_2 \geq 1 \quad and \quad \bar{\tau}_k = \max_{i \in \{0,\ldots,k\}} \tau_i.$$

*Then* $\hat{\alpha}_{sk+j}$, $\hat{\beta}_{sk+j+1}$, *and* $\hat{v}_{sk+j+1}$ *will be computed such that*

$$A\hat{V}_{sk+j} = \hat{V}_{sk+j}\hat{T}_{sk+j} + \hat{\beta}_{sk+j+1}\hat{v}_{sk+j+1}e_{sk+j+1}^T - \delta\hat{V}_{sk+j},$$

*with*

$$\hat{V}_{sk+j} = [\hat{v}_0, \hat{v}_1, \ldots, \hat{v}_{sk+j}]$$

$$\delta\hat{V}_{sk+j} = [\delta\hat{v}_0, \delta\hat{v}_1, \ldots, \delta\hat{v}_{sk+j}]$$

$$\hat{T}_{sk+j} = \begin{bmatrix} \hat{\alpha}_0 & \hat{\beta}_1 & & \\ \hat{\beta}_1 & \ddots & \ddots & \\ & \ddots & \ddots & \hat{\beta}_{sk+j} \\ & & \hat{\beta}_{sk+j} & \hat{\alpha}_{sk+j} \end{bmatrix}$$

*and*

$$\|\delta\hat{v}_{sk+j}\|_2 \leq \epsilon\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (10s+16)\big)\bar{\Gamma}_k\|A\|_2, \qquad (4.1)$$

$$\hat{\beta}_{sk+j+1}|\hat{v}_{sk+j}^T \hat{v}_{sk+j+1}| \leq 2\epsilon(n+11s+15)\|A\|_2\bar{\Gamma}_k^2, \qquad (4.2)$$

$$|\hat{v}_{sk+j+1}^T \hat{v}_{sk+j+1} - 1| \leq \epsilon(n+8s+12)\bar{\Gamma}_k^2, \quad and \qquad (4.3)$$

$$\left|\hat{\beta}_{sk+j+1}^2 + \hat{\alpha}_{sk+j}^2 + \hat{\beta}_{sk+j}^2 - \|A\hat{v}_{sk+j}\|_2^2\right| \leq$$
$$4\epsilon(sk+j+2)\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (3n+40s+58)\big)\bar{\Gamma}_k^2\|A\|_2^2. \quad (4.4)$$

    *Furthermore, if* $R_{sk+j}$ *is the strictly upper triangular matrix such that*

$$\hat{V}_{sk+j}^T \hat{V}_{sk+j} = R_{sk+j}^T + diag(\hat{V}_{sk+j}^T \hat{V}_{sk+j}) + R_{sk+j},$$

*then*

$$\hat{T}_{sk+j}R_{sk+j} - R_{sk+j}\hat{T}_{sk+j} = \hat{\beta}_{sk+j+1}\hat{V}_{sk+j}^T\hat{v}_{sk+j+1}e_{sk+j+1}^T + H_{sk+j}, \qquad (4.5)$$

*where $H_{sk+j}$ is upper triangular with elements $\eta$ such that*

$$
\begin{aligned}
|\eta_{1,1}| &\leq 2\epsilon(n{+}11s{+}15)\|A\|_2\bar{\Gamma}_k^2, \quad and\ for\ i \in \{2,\ldots,sk{+}j{+}1\}, \\
|\eta_{i,i}| &\leq 4\epsilon(n{+}11s{+}15)\|A\|_2\bar{\Gamma}_k^2, \\
|\eta_{i-1,i}| &\leq 2\epsilon\big((n{+}2s{+}5)\sigma{+}(4s{+}9)\bar{\tau}_k + n{+}18s{+}28\big)\bar{\Gamma}_k^2\|A\|_2, \quad and \\
|\eta_{\ell,i}| &\leq 2\epsilon\big((n{+}2s{+}5)\sigma{+}(4s{+}9)\bar{\tau}_k{+}(10s{+}16)\big)\bar{\Gamma}_k^2\|A\|_2, \ for\ \ell \in \{1,\ldots,i{-}2\}.
\end{aligned}
\qquad (4.6)
$$

*Remarks.* This generalizes Paige [27] as follows. The bounds in Theorem 4.2 give insight into how orthogonality is lost in the finite precision $s$-step Lanczos algorithm. Equation (4.1) bounds the error in the columns of the resulting perturbed Lanczos recurrence. How far the Lanczos vectors can deviate from unit 2-norm is given in (4.3), and (4.2) bounds how far adjacent vectors are from being orthogonal. The bound in (4.4) describes how close the columns of $A\hat{V}_{sk+j}$ and $\hat{T}_{sk+j}$ are in size. Finally, (4.5) can be thought of as a recurrence for the loss of orthogonality between Lanczos vectors, and shows how errors propagate through the iterations.

One thing to notice about the bounds in Theorem 4.2 is that they depend heavily on the term $\bar{\Gamma}_k$, which is a measure of the conditioning of the computed $s$-step Krylov bases. This indicates that if $\bar{\Gamma}_k$ is controlled in some way to be near constant, i.e., $\bar{\Gamma}_k = O(1)$, the bounds in Theorem 4.2 will be on the same order as Paige's analogous bounds for classical Lanczos [27], and thus we can expect orthogonality to be lost at a similar rate. The bounds also suggest that for the $s$-step variant to have any use, we must have $\bar{\Gamma}_k = o(\epsilon^{-1/2})$. Otherwise there can be no expectation of orthogonality. Note that $\||\mathcal{B}_k|\|_2$ should be $\lesssim \||A|\|_2$ for many practical basis choices.

Comparing to Paige's result, we can think of $sk + j$ steps of classical Lanczos as the case where $s = 1$, with $\mathcal{Y}_0 = I_{n,n}$ (and then $v_{sk+j} = v'_{k,j}$, $\mathcal{B}_k = A$). In this case $\bar{\Gamma}_k = 1$ and $\bar{\tau}_k = \sigma$ and our bounds reduce (modulo constants) to those of Paige [27].

**4.1. Proof of Theorem 4.2.** The remainder of this section is dedicated to the proof of Theorem 4.2. We first proceed toward proving (4.3).

In finite precision, the Gram matrix construction in line 4 of Algorithm 2 becomes

$$\hat{G}_k = fl(\hat{\mathcal{Y}}_k^T\hat{\mathcal{Y}}_k) = \hat{\mathcal{Y}}_k^T\hat{\mathcal{Y}}_k + \delta G_k, \quad where \quad |\delta G_k| \leq \epsilon n|\hat{\mathcal{Y}}_k^T||\hat{\mathcal{Y}}_k|, \qquad (4.7)$$

and line 14 of Algorithm 2, becomes $\hat{\beta}_{sk+j+1} = fl\big(\ fl(\hat{w}_{k,j}'^T\hat{G}_k\hat{w}'_{k,j})^{1/2}\big)$. Let

$$
\begin{aligned}
d &= fl(\hat{w}_{k,j}'^T\hat{G}_k\hat{w}'_{k,j}) = (\hat{w}_{k,j}'^T + \delta\hat{w}_{k,j}'^T)(\hat{G}_k + \delta\hat{G}_{k,w_j})\hat{w}'_{k,j} \\
&= \hat{w}_{k,j}'^T\hat{\mathcal{Y}}_k^T\hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{w}_{k,j}'^T\delta G_k\hat{w}'_{k,j} + \hat{w}_{k,j}'^T\delta\hat{G}_{k,w_j}\hat{w}'_{k,j} + \delta\hat{w}_{k,j}'^T\hat{G}_k\hat{w}'_{k,j},
\end{aligned}
$$

where

$$|\delta\hat{w}'_{k,j}| \leq \epsilon(2s{+}2)|\hat{w}'_{k,j}| \quad and \qquad\qquad (4.8)$$

$$|\delta\hat{G}_{k,w_j}| \leq \epsilon(2s{+}2)|\hat{G}_k|. \qquad\qquad (4.9)$$

Remember that in the above equation we have ignored all $\epsilon^2$ terms. Now, we let $c = \hat{w}_{k,j}'^T\delta G_k\hat{w}'_{k,j} + \hat{w}_{k,j}'^T\delta\hat{G}_{k,w_j}\hat{w}'_{k,j} + \delta\hat{w}_{k,j}'^T\hat{G}_k\hat{w}'_{k,j}$, where

$$|c| \leq \epsilon(n{+}4s{+}4)\Gamma_k^2\|\hat{\mathcal{Y}}_k\hat{w}'_{k,j}\|_2^2. \qquad\qquad (4.10)$$

We can then write

$$d = \left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2 + c = \left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2 + c \cdot \frac{\left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2}{\left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2} = \left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2 \left(1 + \frac{c}{\left\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\right\|_2^2}\right),$$

and the computation of $\hat{\beta}_{sk+j+1}$ becomes

$$\hat{\beta}_{sk+j+1} = fl(\sqrt{d}) = \sqrt{d} + \delta\beta_{sk+j+1} = \|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2\left(1 + \frac{c}{2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2}\right) + \delta\beta_{sk+j+1}, \quad (4.11)$$

where

$$|\delta\beta_{sk+j+1}| \leq \epsilon\sqrt{d} = \epsilon\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2. \quad (4.12)$$

The coordinate vector $\hat{v}'_{k,j+1}$ is computed as

$$\hat{v}'_{k,j+1} = fl(\hat{w}'_{k,j}/\hat{\beta}_{sk+j+1}) = (\hat{w}'_{k,j} + \delta\tilde{w}'_{k,j})/\hat{\beta}_{sk+j+1}, \quad (4.13)$$

where

$$|\delta\tilde{w}'_{k,j}| \leq \epsilon|\hat{w}'_{k,j}|. \quad (4.14)$$

The corresponding Lanczos vector $\hat{v}_{sk+j+1}$ (as well as $\hat{u}_{sk+j+1}$) are recovered by a change of basis: in finite precision, we have

$$\hat{v}_{sk+j+1} = fl(\hat{\mathcal{Y}}_k \hat{v}'_{k,j+1}) = \left(\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,v_{j+1}}\right)\hat{v}'_{k,j+1}, \quad |\delta\hat{\mathcal{Y}}_{k,v_{j+1}}| \leq \epsilon(2s+2)|\hat{\mathcal{Y}}_k|, \quad (4.15)$$

and

$$\hat{u}_{sk+j+1} = fl(\hat{\mathcal{Y}}_k \hat{u}'_{k,j+1}) = \left(\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,u_{j+1}}\right)\hat{u}'_{k,j+1}, \quad |\delta\hat{\mathcal{Y}}_{k,u_{j+1}}| \leq \epsilon(2s+2)|\hat{\mathcal{Y}}_k|. \quad (4.16)$$

We can now prove (4.3) in Theorem 4.2. Using (4.11), (4.13), and (4.15),

$$\hat{v}_{sk+j+1}^T \hat{v}_{sk+j+1} = \hat{v}_{k,j+1}'^T (\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,v_{j+1}})^T (\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,v_{j+1}})\hat{v}'_{k,j+1}$$

$$= \left(\frac{\hat{w}'_{k,j} + \delta\tilde{w}'_{k,j}}{\hat{\beta}_{sk+j+1}}\right)^T (\hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k + 2\delta\hat{\mathcal{Y}}_{k,v_{j+1}}^T \hat{\mathcal{Y}}_k)\left(\frac{\hat{w}'_{k,j} + \delta\tilde{w}'_{k,j}}{\hat{\beta}_{sk+j+1}}\right)$$

$$= \frac{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + 2\hat{w}_{k,j}'^T \delta\hat{\mathcal{Y}}_{k,v_{j+1}}^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j} + 2\delta\tilde{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j}}{\hat{\beta}_{sk+j+1}^2}$$

$$= \frac{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + 2\hat{w}_{k,j}'^T \delta\hat{\mathcal{Y}}_{k,v_{j+1}}^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j} + 2\delta\tilde{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j}}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + (c + 2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2 \cdot \delta\beta_{sk+j+1})}$$

$$= \frac{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^4}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^4} - \frac{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2(c + 2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2 \cdot \delta\beta_{sk+j+1})}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^4}$$

$$+ \frac{2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2(\hat{w}_{k,j}'^T \delta\hat{\mathcal{Y}}_{k,v_{j+1}}^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j} + \delta\tilde{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j})}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^4}$$

$$= 1 - \frac{c + 2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2 \cdot \delta\beta_{sk+j+1}}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2}$$

$$+ \frac{2(\hat{w}_{k,j}'^T \delta\hat{\mathcal{Y}}_{k,v_{j+1}}^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j} + \delta\tilde{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}'_{k,j})}{\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2}$$

Now, using bounds in (4.7), (4.8), (4.9), (4.10), (4.15), (4.16), and Lemma 4.1, we obtain

$$\begin{aligned}
|\hat{v}_{sk+j+1}^T \hat{v}_{sk+j+1} - 1| &\leq \epsilon(n{+}4s{+}4)\Gamma_k^2 + 2\epsilon + 2\epsilon(2s{+}2)\Gamma_k + 2\epsilon\Gamma_k \\
&\leq \epsilon(n{+}4s{+}4)\Gamma_k^2 + \epsilon(4s{+}6)\Gamma_k + 2\epsilon \\
&\leq \epsilon(n{+}8s{+}12)\Gamma_k^2.
\end{aligned}$$

This thus proves (4.3), and we now proceed toward proving (4.2). Using (4.9), line 12 in Algorithm 2 is computed in finite precision as

$$\hat{\alpha}_{sk+j} = fl(\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}') = (\hat{v}_{k,j}'^T + \delta\hat{v}_{k,j}'^T)(\hat{G}_k + \delta\hat{G}_{k,u_j})\hat{u}_{k,j}',$$

where $|\delta\hat{v}_{k,j}'| \leq \epsilon(2s{+}2)|\hat{v}_{k,j}'|$ and $|\delta\hat{G}_{k,u_j}| \leq \epsilon(2s + 2)|\hat{G}_k|$. Expanding the above equation using (4.7), and (4.15), we obtain

$$\begin{aligned}
\hat{\alpha}_{sk+j} =&\, \hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' + \hat{v}_{k,j}'^T \delta\hat{G}_{k,u_j} \hat{u}_{k,j}' + \delta\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' \\
=&\, \hat{v}_{k,j}'^T (\hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k + \delta G_k)\hat{u}_{k,j}' + \hat{v}_{k,j}'^T \delta\hat{G}_{k,u_j}\hat{u}_{k,j}' + \delta\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' \\
=&\, \hat{v}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{u}_{k,j}' + \hat{v}_{k,j}'^T \delta G_k \hat{u}_{k,j}' + \hat{v}_{k,j}'^T \delta\hat{G}_{k,u_j}\hat{u}_{k,j}' + \delta\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' \\
=&\, (\hat{v}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}_{k,j}')^T(\hat{u}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,u_j}\hat{u}_{k,j}') + \hat{v}_{k,j}'^T \delta G_k \hat{u}_{k,j}' + \hat{v}_{k,j}'^T \delta\hat{G}_{k,u_j}\hat{u}_{k,j}' \\
&+ \delta\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' \\
=&\, \hat{v}_{sk+j}^T \hat{u}_{sk+j} + \delta\hat{\alpha}_{sk+j}, \tag{4.17}
\end{aligned}$$

with $\delta\hat{\alpha}_{sk+j} = \delta\hat{v}_{k,j}'^T \hat{G}_k \hat{u}_{k,j}' + \hat{v}_{k,j}'^T(\delta G_k + \delta\hat{G}_{k,u_j} - \hat{\mathcal{Y}}_k^T \delta\hat{\mathcal{Y}}_{k,u_j} - \delta\hat{\mathcal{Y}}_{k,v_j}^T \hat{\mathcal{Y}}_k)\hat{u}_{k,j}'$.

Using bounds in (4.3), (4.7), (4.8), (4.9), (4.15), and (4.16), as well as Lemma 4.1, we can write (again, ignoring $\epsilon^2$ terms),

$$\begin{aligned}
|\delta\hat{\alpha}_{sk+j}| &\leq \epsilon(n{+}8s{+}8)|\hat{v}_{k,j}'^T||\hat{\mathcal{Y}}_k^T||\hat{\mathcal{Y}}_k||\hat{u}_{k,j}'| \\
&\leq \epsilon(n{+}8s{+}8) \,\||\hat{\mathcal{Y}}_k||\hat{v}_{k,j}'|\|_2 \,\||\hat{\mathcal{Y}}_k||\hat{u}_{k,j}'|\|_2 \\
&\leq \epsilon(n{+}8s{+}8)(\Gamma_k\|\hat{v}_{sk+j}\|_2)(\Gamma_k\|\hat{u}_{sk+j}\|_2) \\
&\leq \epsilon(n{+}8s{+}8)\big(\Gamma_k(1 + (\epsilon/2)(n{+}8s{+}12)\Gamma_k^2)\big)(\Gamma_k\|\hat{u}_{sk+j}\|_2) \\
&\leq \epsilon(n{+}8s{+}8)\Gamma_k(\Gamma_k\|\hat{u}_{sk+j}\|_2) \\
&\leq \epsilon(n{+}8s{+}8)\Gamma_k^2\|\hat{u}_{sk+j}\|_2. \tag{4.18}
\end{aligned}$$

Taking the norm of (4.17), and using the bounds in (4.18) and (4.3), we obtain the bound

$$\begin{aligned}
|\hat{\alpha}_{sk+j}| &\leq \|\hat{v}_{sk+j}^T\|_2\|\hat{u}_{sk+j}\|_2 + |\delta\hat{\alpha}_{sk+j}| \\
&\leq \big(1 + (\epsilon/2)(n{+}8s{+}12)\Gamma_k^2\big)\|\hat{u}_{sk+j}\|_2 + \epsilon(n{+}8s{+}8)\Gamma_k^2\|\hat{u}_{sk+j}\|_2 \\
&\leq \big(1 + \epsilon\big((3/2)n{+}12s{+}14\big)\Gamma_k^2\big)\|\hat{u}_{sk+j}\|_2. \tag{4.19}
\end{aligned}$$

In finite precision, line 13 of Algorithm 2 is computed as

$$\hat{w}_{k,j}' = \hat{u}_{k,j}' - \hat{\alpha}_{sk+j}\hat{v}_{k,j}' - \delta w_{k,j}', \quad \text{where} \quad |\delta w_{k,j}'| \leq \epsilon(|\hat{u}_{k,j}'| + 2|\hat{\alpha}_{sk+j}\hat{v}_{k,j}'|). \tag{4.20}$$

Multiplying both sides of (4.20) by $\hat{\mathcal{Y}}_k$ gives

$$\hat{\mathcal{Y}}_k\hat{w}_{k,j}' = \hat{\mathcal{Y}}_k\hat{u}_{k,j}' - \hat{\alpha}_{sk+j}\hat{\mathcal{Y}}_k\hat{v}_{k,j}' - \hat{\mathcal{Y}}_k\delta w_{k,j}',$$

and multiplying each side by its own transpose, we get

$$
\begin{aligned}
\hat{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}' &= (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}' - \hat{\mathcal{Y}}_k \delta w_{k,j}')^T (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}' - \hat{\mathcal{Y}}_k \delta w_{k,j}') \\
&= \hat{u}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{u}_{k,j}' - 2\hat{\alpha}_{sk+j} \hat{u}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{v}_{k,j}' + \hat{\alpha}_{sk+j}^2 \hat{v}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{v}_{k,j}' \\
&\quad - \delta w_{k,j}'^T \hat{\mathcal{Y}}_k^T (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}') - (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}')^T \hat{\mathcal{Y}}_k \delta w_{k,j}'.
\end{aligned}
$$

Using (4.15) and (4.16), we can write

$$
\begin{aligned}
\hat{w}_{k,j}'^T \hat{\mathcal{Y}}_k^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}' &= (\hat{u}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}')^T (\hat{u}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}') \\
&\quad - 2\hat{\alpha}_{sk+j} (\hat{u}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}')^T (\hat{v}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}') \\
&\quad + \hat{\alpha}_{sk+j}^2 (\hat{v}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}')^T (\hat{v}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}') \\
&\quad - 2\delta w_{k,j}'^T \hat{\mathcal{Y}}_k^T (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}') \\
&= \hat{u}_{sk+j}^T \hat{u}_{sk+j} - 2\hat{u}_{sk+j}^T \delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - 2\hat{\alpha}_{sk+j} \hat{u}_{sk+j}^T \hat{v}_{sk+j} \\
&\quad + 2\hat{\alpha}_{sk+j} \hat{u}_{sk+j}^T \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' + 2\hat{\alpha}_{sk+j} \hat{u}_{k,j}'^T \delta \hat{\mathcal{Y}}_{k,u_j}^T \hat{v}_{sk+j} \\
&\quad + \hat{\alpha}_{sk+j}^2 \hat{v}_{sk+j}^T \hat{v}_{sk+j} - 2\hat{\alpha}_{sk+j}^2 \hat{v}_{sk+j}^T \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' \\
&\quad - 2\delta w_{k,j}'^T \hat{\mathcal{Y}}_k^T (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}') \\
&= \hat{u}_{sk+j}^T \hat{u}_{sk+j} - 2\hat{\alpha}_{sk+j} \hat{u}_{sk+j}^T \hat{v}_{sk+j} + \hat{\alpha}_{sk+j}^2 \hat{v}_{sk+j}^T \hat{v}_{sk+j} \\
&\quad - 2(\delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}')^T (\hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j}) \\
&\quad - 2\delta w_{k,j}'^T \hat{\mathcal{Y}}_k^T (\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}').
\end{aligned}
$$

This can be written

$$
\begin{aligned}
\|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 &= \|\hat{u}_{sk+j}\|_2^2 - 2\hat{\alpha}_{sk+j} \hat{u}_{sk+j}^T \hat{v}_{sk+j} + \hat{\alpha}_{sk+j}^2 \|\hat{v}_{sk+j}\|_2^2 \\
&\quad - 2(\delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' + \hat{\mathcal{Y}}_k \delta w_{k,j}')^T (\hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j}),
\end{aligned}
$$

where we have used $\hat{\mathcal{Y}}_k \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \hat{\mathcal{Y}}_k \hat{v}_{k,j}' = \hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j} + O(\epsilon)$. Now, using (4.17),

$$
\begin{aligned}
\|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 &= \|\hat{u}_{sk+j}\|_2^2 - 2\hat{\alpha}_{sk+j} (\hat{\alpha}_{sk+j} - \delta \hat{\alpha}_{sk+j}) + \hat{\alpha}_{sk+j}^2 \|\hat{v}_{sk+j}\|_2^2 \\
&\quad - 2(\delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' + \hat{\mathcal{Y}}_k \delta w_{k,j}')^T (\hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j}) \\
&= \|\hat{u}_{sk+j}\|_2^2 + \hat{\alpha}_{sk+j}^2 (\|\hat{v}_{sk+j}\|_2^2 - 2) + 2\hat{\alpha}_{sk+j} \delta \hat{\alpha}_{sk+j} \\
&\quad - 2(\delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' + \hat{\mathcal{Y}}_k \delta w_{k,j}')^T (\hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j}).
\end{aligned}
$$

Now, we rearrange the above equation to obtain

$$
\begin{aligned}
\|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 &= \hat{\alpha}_{sk+j}^2 (\|\hat{v}_{sk+j}\|_2^2 - 1) + 2\hat{\alpha}_{sk+j} \delta \hat{\alpha}_{sk+j} \\
&\quad - 2(\delta \hat{\mathcal{Y}}_{k,u_j} \hat{u}_{k,j}' - \hat{\alpha}_{sk+j} \delta \hat{\mathcal{Y}}_{k,v_j} \hat{v}_{k,j}' + \hat{\mathcal{Y}}_k \delta w_{k,j}')^T (\hat{u}_{sk+j} - \hat{\alpha}_{sk+j} \hat{v}_{sk+j}).
\end{aligned}
$$

Using Lemma 4.1 and bounds in (4.3), (4.15), (4.16), (4.18), (4.19), and (4.20),

we can then write

$$\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 \leq \left(1 + \epsilon\left((3/2)n + 12s + 14\right)\Gamma_k^2\right)^2 \cdot \epsilon(n + 8s + 12)\Gamma_k^2\|\hat{u}_{sk+j}\|_2^2$$
$$+ 2\left(1 + \epsilon\left((3/2)n + 12s + 14\right)\Gamma_k^2\right) \cdot \epsilon(n + 8s + 8)\Gamma_k^2\|\hat{u}_{sk+j}\|_2^2$$
$$+ 2\epsilon\left((2s+2) + (2s+2) + 3\right)\Gamma_k\|\hat{u}_{sk+j}\|_2 \cdot 2\|\hat{u}_{sk+j}\|_2$$
$$\leq \epsilon(n+8s+12)\Gamma_k^2\|\hat{u}_{sk+j}\|_2^2 + 2\epsilon(n+8s+8)\Gamma_k^2\|\hat{u}_{sk+j}\|_2^2$$
$$+ \epsilon(16s+28)\Gamma_k\|\hat{u}_{sk+j}\|_2^2,$$

where, again, we have ignored $\epsilon^2$ terms. Using $\Gamma_k \leq \Gamma_k^2$, this gives the bound

$$\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 \leq \epsilon(3n+40s+56)\Gamma_k^2\|\hat{u}_{sk+j}\|_2^2. \qquad (4.21)$$

Given the above, we can also write the bound

$$\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 \leq \|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2 + \hat{\alpha}_{sk+j}^2 \leq \left(1 + \epsilon(3n+40s+56)\Gamma_k^2\right)\|\hat{u}_{sk+j}\|_2^2, \qquad (4.22)$$

and using (4.10), (4.11), and (4.12),

$$|\hat{\beta}_{sk+j+1}| \leq \|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2 \left(1 + \epsilon + \frac{c}{2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2}\right)$$
$$\leq \left(1 + (1/2)\epsilon(3n+40s+56)\Gamma_k^2\right)\|\hat{u}_{sk+j}\|_2 \left(1 + \epsilon + \frac{c}{2\|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2^2}\right)$$
$$\leq \left(1 + \epsilon + (1/2)\epsilon(n+4s+4)\Gamma_k^2 + (1/2)\epsilon(3n+40s+56)\Gamma_k^2\right)\|\hat{u}_{sk+j}\|_2.$$

Combining terms and using $1 \leq \Gamma_k^2$, the above can be written

$$|\hat{\beta}_{sk+j+1}| \leq \left(1 + \epsilon(2n+22s+31)\Gamma_k^2\right)\|\hat{u}_{sk+j}\|_2. \qquad (4.23)$$

Now, rearranging (4.13), we can write

$$\hat{\beta}_{sk+j+1}\hat{v}'_{k,j+1} = \hat{w}'_{k,j} + \delta\tilde{w}'_{k,j},$$

and premultiplying by $\hat{\mathcal{Y}}_k$, we obtain

$$\hat{\beta}_{sk+j+1}\hat{\mathcal{Y}}_k\hat{v}'_{k,j+1} = \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{\mathcal{Y}}_k\delta\tilde{w}'_{k,j}.$$

Using (4.15), this can be written

$$\hat{\beta}_{sk+j+1}(\hat{v}_{sk+j+1} - \delta\hat{\mathcal{Y}}_{k,v_{j+1}}\hat{v}'_{k,j+1}) = \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{\mathcal{Y}}_k\delta\tilde{w}'_{k,j}.$$

Rearranging and using (4.13),

$$\hat{\beta}_{sk+j+1}\hat{v}_{sk+j+1} = \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{\mathcal{Y}}_k\delta\tilde{w}'_{k,j} + \hat{\beta}_{sk+j+1}\delta\hat{\mathcal{Y}}_{k,v_{j+1}}\hat{v}'_{k,j+1}$$
$$= \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{\mathcal{Y}}_k\delta\tilde{w}'_{k,j} + \delta\hat{\mathcal{Y}}_{k,v_{j+1}}(\hat{w}'_{k,j} + \delta\tilde{w}'_{k,j})$$
$$= \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \hat{\mathcal{Y}}_k\delta\tilde{w}'_{k,j} + \delta\hat{\mathcal{Y}}_{k,v_{j+1}}\hat{w}'_{k,j}$$
$$\equiv \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \delta w_{sk+j}, \qquad (4.24)$$

where $\delta w_{sk+j} = \hat{\mathcal{Y}}_k \delta \tilde{w}'_{k,j} + \delta \hat{\mathcal{Y}}_{k,v_{j+1}} \hat{w}'_{k,j}$. Using Lemma 4.1 and bounds in (4.14), (4.15), and (4.22),

$$
\begin{aligned}
\|\delta w_{sk+j}\|_2 &\le \epsilon \| |\hat{\mathcal{Y}}_k| |\hat{w}'_{k,j}| \|_2 + \epsilon(2s+2) \| |\hat{\mathcal{Y}}_k| |\hat{w}'_{k,j}| \|_2 \\
&\le \epsilon(2s+3)\Gamma_k \|\hat{\mathcal{Y}}_k \hat{w}'_{k,j}\|_2 \\
&\le \epsilon(2s+3)\Gamma_k \|\hat{u}_{sk+j}\|_2.
\end{aligned}
\tag{4.25}
$$

We premultiply (4.24) by $\hat{v}^T_{sk+j}$ and use (4.15), (4.16), (4.17), and (4.20) to obtain

$$
\begin{aligned}
\hat{\beta}_{sk+j+1}\hat{v}^T_{sk+j}\hat{v}_{sk+j+1} &= \hat{v}^T_{sk+j}(\hat{\mathcal{Y}}_k \hat{w}'_{k,j} + \delta w_{sk+j}) \\
&= \hat{v}^T_{sk+j}(\hat{\mathcal{Y}}_k \hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\hat{\mathcal{Y}}_k \hat{v}'_{k,j} - \hat{\mathcal{Y}}_k \delta \hat{w}'_{k,j}) + \hat{v}^T_{sk+j}\delta w_{sk+j} \\
&= \hat{v}^T_{sk+j}(\hat{\mathcal{Y}}_k \hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\hat{\mathcal{Y}}_k \hat{v}'_{k,j}) - \hat{v}^T_{sk+j}(\hat{\mathcal{Y}}_k \delta w'_{k,j} - \delta w_{sk+j}) \\
&= \hat{v}^T_{sk+j}((\hat{u}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j}) - \hat{\alpha}_{sk+j}(\hat{v}_{sk+j} - \delta \hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j})) \\
&\quad - \hat{v}^T_{sk+j}(\hat{\mathcal{Y}}_k \delta w'_{k,j} - \delta w_{sk+j}) \\
&= \hat{v}^T_{sk+j}\hat{u}_{sk+j} - \hat{\alpha}_{sk+j}\hat{v}^T_{sk+j}\hat{v}_{sk+j} \\
&\quad - \hat{v}^T_{sk+j}(\delta \hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\delta \hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} + \hat{\mathcal{Y}}_k \delta w'_{k,j} - \delta w_{sk+j}) \\
&= (\hat{\alpha}_{sk+j} - \delta \hat{\alpha}_{sk+j}) - \hat{\alpha}_{sk+j}\|\hat{v}_{sk+j}\|_2^2 \\
&\quad - \hat{v}^T_{sk+j}(\delta \hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\delta \hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} + \hat{\mathcal{Y}}_k \delta w'_{k,j} - \delta w_{sk+j}) \\
&= -\delta \hat{\alpha}_{sk+j} - \hat{\alpha}_{sk+j}(\|\hat{v}_{sk+j}\|_2^2 - 1) \\
&\quad - \hat{v}^T_{sk+j}(\delta \hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\delta \hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} + \hat{\mathcal{Y}}_k \delta w'_{k,j} - \delta w_{sk+j}),
\end{aligned}
$$

and using Lemma 4.1 and bounds in (4.3), (4.13), (4.15), (4.16), (4.18), (4.19), (4.20), and (4.25), we can write the bound

$$
\begin{aligned}
\left| \hat{\beta}_{sk+j+1} \cdot \hat{v}^T_{sk+j}\hat{v}_{sk+j+1} \right| &\le |\delta \hat{\alpha}_{sk+j}| + |\hat{\alpha}_{sk+j}||\hat{v}^T_{sk+j}\hat{v}_{sk+j} - 1| \\
&\quad + \|\hat{v}_{sk+j}\|_2 (\| |\delta \hat{\mathcal{Y}}_{k,u_j}| |\hat{u}'_{k,j}| \|_2 + |\hat{\alpha}_{sk+j}| \| |\delta \hat{\mathcal{Y}}_{k,v_j}| |\hat{v}'_{k,j}| \|_2) \\
&\quad + \|\hat{v}_{sk+j}\|_2 (\| |\hat{\mathcal{Y}}_k| |\delta w'_{k,j}| \|_2 + \|\delta w_{sk+j}\|_2) \\
&\le 2\epsilon(n+11s+15)\Gamma_k^2 \|\hat{u}_{sk+j}\|_2.
\end{aligned}
\tag{4.26}
$$

This is a start toward proving (4.2). We will return to the above bound once we later prove a bound on $\|\hat{u}_{sk+j}\|_2$. Our next step is to analyze the error in each column of the finite precision $s$-step Lanczos recurrence. First, we note that we can write the error in computing the $s$-step bases (line 3 in Algorithm 2) by

$$
A\underline{\hat{\mathcal{Y}}}_k = \hat{\mathcal{Y}}_k \mathcal{B}_k + \delta E_k
\tag{4.27}
$$

where $\underline{\hat{\mathcal{Y}}}_k = [\hat{\mathcal{V}}_k[I_s, 0_{s,1}]^T, 0_{n,1}, \hat{\mathcal{U}}_k[I_s, 0_{s,1}]^T, 0_{n,1}]$. It can be shown (see, e.g., [3]) that if the basis is computed in the usual way by repeated SpMVs,

$$
|\delta E_k| \le \epsilon((3+n)|A||\underline{\hat{\mathcal{Y}}}_k| + 4|\hat{\mathcal{Y}}_k||\mathcal{B}_k|).
\tag{4.28}
$$

In finite precision, line 17 in Algorithm 2 is computed as

$$
\hat{u}'_{k,j} = \mathcal{B}_k \hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{v}'_{k,j-1} + \delta u'_{k,j}, \quad |\delta u'_{k,j}| \le \epsilon((2s+3)|\mathcal{B}_k||\hat{v}'_{k,j}| + 2|\hat{\beta}_{sk+j}\hat{v}'_{k,j-1}|),
\tag{4.29}
$$

and then, with Lemma 4.1, (4.15), (4.16), (4.27), and (4.29), we can write

$$
\begin{aligned}
\hat{u}_{sk+j} &= (\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,u_j})\hat{u}'_{k,j} \\
&= (\hat{\mathcal{Y}}_k + \delta\hat{\mathcal{Y}}_{k,u_j})(\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{v}'_{k,j-1} + \delta u'_{k,j}) \\
&= \hat{\mathcal{Y}}_k\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{\mathcal{Y}}_k\hat{v}'_{k,j-1} + \hat{\mathcal{Y}}_k\delta u'_{k,j} + \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,u_j}\hat{v}'_{k,j-1} \\
&= (A\underline{\hat{\mathcal{Y}}}_k - \delta E_k)\hat{v}'_{k,j} - \hat{\beta}_{sk+j}(\hat{v}_{sk+j-1} - \delta\hat{\mathcal{Y}}_{k,v_{j-1}}\hat{v}'_{k,j-1}) + \hat{\mathcal{Y}}_k\delta u'_{k,j} \\
&\quad + \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,u_j}\hat{v}'_{k,j-1} \\
&= A\underline{\hat{\mathcal{Y}}}_k\hat{v}'_{k,j} - \delta E_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1} + \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,v_{j-1}}\hat{v}'_{k,j-1} + \hat{\mathcal{Y}}_k\delta u'_{k,j} \\
&\quad + \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,u_j}\hat{v}'_{k,j-1} \\
&= A(\hat{v}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j}) - \delta E_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1} + \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,v_{j-1}}\hat{v}'_{k,j-1} \\
&\quad + \hat{\mathcal{Y}}_k\delta u'_{k,j} + \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,u_j}\hat{v}'_{k,j-1} \\
&= A\hat{v}_{sk+j} - A\delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} - \delta E_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1} + \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,v_{j-1}}\hat{v}'_{k,j-1} \\
&\quad + \hat{\mathcal{Y}}_k\delta u'_{k,j} + \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k\hat{v}'_{k,j} - \hat{\beta}_{sk+j}\delta\hat{\mathcal{Y}}_{k,u_j}\hat{v}'_{k,j-1} \\
&\equiv A\hat{v}_{sk+j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1} + \delta u_{sk+j}, \quad\quad\quad\quad (4.30)
\end{aligned}
$$

where

$$
\delta u_{sk+j} = \hat{\mathcal{Y}}_k\delta u'_{k,j} - (A\delta\hat{\mathcal{Y}}_{k,v_j} - \delta\hat{\mathcal{Y}}_{k,u_j}\mathcal{B}_k + \delta E_k)\hat{v}'_{k,j} + \hat{\beta}_{sk+j}(\delta\hat{\mathcal{Y}}_{k,v_{j-1}} - \delta\hat{\mathcal{Y}}_{k,u_j})\hat{v}'_{k,j-1}.
$$

Using the bounds in (4.15), (4.16), (4.23), (4.28), and (4.29) we can write

$$
\begin{aligned}
|\delta u_{sk+j}| &\leq \epsilon\big((2s+3)|\hat{\mathcal{Y}}_k||\mathcal{B}_k||\hat{v}'_{k,j}| + 2|\hat{\beta}_{sk+j}||\hat{\mathcal{Y}}_k||\hat{v}'_{k,j-1}|\big) \\
&\quad + \epsilon(2s+2)|A||\hat{\mathcal{Y}}_k||\hat{v}'_{k,j}| + \epsilon(2s+2)|\hat{\mathcal{Y}}_k||\mathcal{B}_k||\hat{v}'_{k,j}| \\
&\quad + \epsilon\big((3+n)|A||\hat{\mathcal{Y}}_k||\hat{v}'_{k,j}| + 4|\hat{\mathcal{Y}}_k||\mathcal{B}_k||\hat{v}'_{k,j}|\big) \\
&\quad + 2\epsilon(2s+2)|\hat{\beta}_{sk+j}||\hat{\mathcal{Y}}_k||\hat{v}'_{k,j-1}| \\
&\leq \epsilon(n+2s+5)|A||\hat{\mathcal{Y}}_k||\hat{v}'_{k,j}| + \epsilon(4s+9)|\hat{\mathcal{Y}}_k||\mathcal{B}_k||\hat{v}'_{k,j}| \\
&\quad + \epsilon(4s+6)\big(1 + \epsilon(2n+22s+31)\Gamma_k^2\big)\|\hat{u}_{sk+j-1}\|_2 \cdot |\hat{\mathcal{Y}}_k||\hat{v}'_{k,j-1}|.
\end{aligned}
$$

and from this we obtain

$$
\begin{aligned}
\|\delta u_{sk+j}\|_2 &\leq \epsilon(n+2s+5)\,\||A|\|_2\,\||\hat{\mathcal{Y}}_k|\|_2\,\|\hat{v}'_{k,j}\|_2 + \epsilon(4s+9)\,\||\hat{\mathcal{Y}}_k|\|_2\,\||\mathcal{B}_k|\|_2\,\|\hat{v}'_{k,j}\|_2 \\
&\quad + \epsilon(4s+6)\,\||\hat{\mathcal{Y}}_k|\|_2\,\|\hat{v}'_{k,j-1}\|_2\,\|\hat{u}_{sk+j-1}\|_2 \\
&\leq \epsilon(n+2s+5)\,\||A|\|_2\,\Gamma_k + \epsilon(4s+9)\,\||\mathcal{B}_k|\|_2\,\Gamma_k + \epsilon(4s+6)\,\|\hat{u}_{sk+j-1}\|_2\,\Gamma_k.
\end{aligned}
$$

We will now introduce and make use of the quantities $\sigma \equiv \||A|\|_2/\|A\|_2$ and $\tau_k \equiv \||\mathcal{B}_k|\|_2/\|A\|_2$. Note that the quantity $\||\mathcal{B}_k|\|_2$ is controlled by the user, and for many popular basis choices, such as monomial, Newton, or Chebyshev bases, it should be the case that $\||\mathcal{B}_k|\|_2 \lesssim \|A\|_2$. Using these quantities, the bound above can be written

$$
\big\|\delta u_{sk+j}\big\|_2 \leq \epsilon\Big(\big((n+2s+5)\sigma + (4s+9)\tau_k\big)\|A\|_2 + (4s+6)\|\hat{u}_{sk+j-1}\|_2\Big)\Gamma_k. \quad (4.31)
$$

Manipulating (4.24), and using (4.15), (4.16), and (4.20), we have

$$
\begin{aligned}
\hat{\beta}_{sk+j+1}\hat{v}_{sk+j+1} &= \hat{\mathcal{Y}}_k\hat{w}'_{k,j} + \delta w_{sk+j} \\
&= \hat{\mathcal{Y}}_k\hat{u}'_{k,j} - \hat{\alpha}_{sk+j}\hat{\mathcal{Y}}_k\hat{v}'_{k,j} - \hat{\mathcal{Y}}_k\delta w'_{k,j} + \delta w_{sk+j} \\
&= (\hat{u}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j}) - \hat{\alpha}_{sk+j}(\hat{v}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j}) - \hat{\mathcal{Y}}_k\delta w'_{k,j} + \delta w_{sk+j} \\
&= \hat{u}_{sk+j} - \hat{\alpha}_{sk+j}\hat{v}_{sk+j} - \delta\hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j} + \hat{\alpha}_{sk+j}\delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} - \hat{\mathcal{Y}}_k\delta w'_{k,j} \\
&\quad + \delta w_{sk+j},
\end{aligned}
$$

and substituting in the expression for $\hat{u}_{sk+j}$ in (4.30) on the right, we obtain

$$
\hat{\beta}_{sk+j+1}\hat{v}_{sk+j+1} \equiv A\hat{v}_{sk+j} - \hat{\alpha}_{sk+j}\hat{v}_{sk+j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1} + \delta\hat{v}_{sk+j}, \qquad (4.32)
$$

where

$$
\delta\hat{v}_{sk+j} = \delta u_{sk+j} - \delta\hat{\mathcal{Y}}_{k,u_j}\hat{u}'_{k,j} + \hat{\alpha}_{sk+j}\delta\hat{\mathcal{Y}}_{k,v_j}\hat{v}'_{k,j} - \hat{\mathcal{Y}}_k\delta w'_{k,j} + \delta w_{sk+j}.
$$

From this we can write the componentwise bound

$$
|\delta\hat{v}_{sk+j}| \leq |\delta u_{sk+j}| + |\delta\hat{\mathcal{Y}}_{k,u_j}|\,|\hat{u}'_{k,j}| + |\hat{\alpha}_{sk+j}|\,|\delta\hat{\mathcal{Y}}_{k,v_j}|\,|\hat{v}'_{k,j}| + |\hat{\mathcal{Y}}_k|\,|\delta w'_{k,j}| + |\delta w_{sk+j}|,
$$

and using Lemma 4.1, (4.14), (4.15), (4.16), (4.19), (4.20), and (4.25) we obtain

$$
\begin{aligned}
\|\delta\hat{v}_{sk+j}\|_2 &\leq \|\delta u_{sk+j}\|_2 + \epsilon(2s+2)\|\,|\hat{\mathcal{Y}}_k||\hat{u}'_{k,j}|\,\|_2 \\
&\quad + \big(1+\epsilon((3/2)n+12s+14)\Gamma_k^2\big)\|\hat{u}_{sk+j}\|_2 \cdot \epsilon(2s+2)\|\,|\hat{\mathcal{Y}}_k||\hat{v}'_{k,j}|\,\|_2 \\
&\quad + \epsilon\|\,|\hat{\mathcal{Y}}_k||\hat{u}'_{k,j}|\,\|_2 + 2\epsilon\big(1+\epsilon((3/2)n+12s+14)\Gamma_k^2\big)\|\hat{u}_{sk+j}\|_2\|\,|\hat{\mathcal{Y}}_k||\hat{v}'_{k,j}|\,\|_2 \\
&\quad + \epsilon(2s+3)\Gamma_k\|\hat{u}_{sk+j}\|_2 \\
&\leq \|\delta u_{sk+j}\|_2 + \epsilon(2s+2)\Gamma_k\|\hat{u}_{sk+j}\|_2 + \epsilon(2s+2)\Gamma_k\|\hat{u}_{sk+j}\|_2 \\
&\quad + \epsilon\Gamma_k\|\hat{u}_{sk+j}\|_2 + 2\epsilon\Gamma_k\|\hat{u}_{sk+j}\|_2 + \epsilon(2s+3)\Gamma_k\|\hat{u}_{sk+j}\|_2 \\
&\leq \|\delta u_{sk+j}\|_2 + \epsilon(6s+10)\Gamma_k\|\hat{u}_{sk+j}\|_2.
\end{aligned}
$$

Using (4.31), this gives the bound

$$
\|\delta\hat{v}_{sk+j}\|_2 \leq
$$
$$
\epsilon\Big(\big((n+2s+5)\sigma+(4s+9)\tau_k\big)\|A\|_2+(6s+10)\|\hat{u}_{sk+j}\|_2+(4s+6)\|\hat{u}_{sk+j-1}\|_2\Big)\Gamma_k. \quad (4.33)
$$

We now have everything we need to write the finite-precision $s$-step Lanczos recurrence in its familiar matrix form. Let

$$
\hat{T}_{sk+j} =
\begin{bmatrix}
\hat{\alpha}_0 & \hat{\beta}_1 & & \\
\hat{\beta}_1 & \ddots & \ddots & \\
& \ddots & \ddots & \hat{\beta}_{sk+j} \\
& & \hat{\beta}_{sk+j} & \hat{\alpha}_{sk+j}
\end{bmatrix},
$$

and let $\hat{V}_{sk+j} = [\hat{v}_0, \hat{v}_1, \ldots, \hat{v}_{sk+j}]$ and $\delta\hat{V}_{sk+j} = [\delta\hat{v}_0, \delta\hat{v}_1, \ldots, \delta\hat{v}_{sk+j}]$. Note that $\hat{T}_{sk+j}$ has dimension $(sk+j+1)$-by-$(sk+j+1)$, and $\hat{V}_{sk+j}$ and $\delta\hat{V}_{sk+j}$ have dimension $n$-by-$(sk+j+1)$. Then (4.32) in matrix form gives

$$
A\hat{V}_{sk+j} = \hat{V}_{sk+j}\hat{T}_{sk+j} + \hat{\beta}_{sk+j+1}\hat{v}_{sk+j+1}e_{sk+j+1}^T - \delta\hat{V}_{sk+j}. \qquad (4.34)
$$

Thus (4.33) gives a bound on the error in the columns of the finite precision $s$-step Lanczos recurrence. Again, we will return to (4.33) to prove (4.1) once we bound $\|\hat{u}_{sk+j}\|_2$.

Now, we examine the possible loss of orthogonality in the vectors $\hat{v}_0, \ldots, \hat{v}_{sk+j+1}$. We define the strictly upper triangular matrix $R_{sk+j}$ of dimension $(sk+j+1)$-by-$(sk+j+1)$ with elements $\rho_{i,j}$, for $i,j \in \{1, \ldots, sk+j+1\}$, such that

$$\hat{V}_{sk+j}^T \hat{V}_{sk+j} = R_{sk+j}^T + \operatorname{diag}(\hat{V}_{sk+j}^T \hat{V}_{sk+j}) + R_{sk+j}.$$

For notational purposes, we also define $\rho_{sk+j+1,sk+j+2} \equiv \hat{v}_{sk+j}^T \hat{v}_{sk+j+1}$. (Note that $\rho_{sk+j+1,sk+j+2}$ is not an element in $R_{sk+j}$, but would be an element in $R_{sk+j+1}$). Multiplying (4.34) on the left by $\hat{V}_{sk+j}^T$, we get

$$\hat{V}_{sk+j}^T A \hat{V}_{sk+j} = \hat{V}_{sk+j}^T \hat{V}_{sk+j} \hat{T}_{sk+j} + \hat{\beta}_{sk+j+1} \hat{V}_{sk+j}^T \hat{v}_{sk+j+1} e_{sk+j+1}^T - \hat{V}_{sk+j}^T \delta\hat{V}_{sk+j}.$$

Since the above is symmetric, we can equate the right hand side by its own transpose to obtain

$$\hat{T}_{sk+j}(R_{sk+j}^T + R_{sk+j}) - (R_{sk+j}^T + R_{sk+j})\hat{T}_{sk+j} =$$
$$\hat{\beta}_{sk+j+1}(\hat{V}_{sk+j}^T \hat{v}_{sk+j+1} e_{sk+j+1}^T - e_{sk+j+1} \hat{v}_{sk+j+1}^T \hat{V}_{sk+j})$$
$$+ \hat{V}_{sk+j}^T \delta\hat{V}_{sk+j} - \delta\hat{V}_{sk+j}^T \hat{V}_{sk+j} + \operatorname{diag}(\hat{V}_{sk+j}^T \hat{V}_{sk+j}) \cdot \hat{T}_{sk+j} - \hat{T}_{sk+j} \cdot \operatorname{diag}(\hat{V}_{sk+j}^T \hat{V}_{sk+j}).$$

Now, let $M_{sk+j} \equiv \hat{T}_{sk+j} R_{sk+j} - R_{sk+j} \hat{T}_{sk+j}$, which is upper triangular and has dimension $(sk+j+1)$-by-$(sk+j+1)$. Then the left-hand side above can be written as $M_{sk+j} - M_{sk+j}^T$, and we can equate the strictly upper triangular part of $M_{sk+j}$ with the strictly upper triangular part of the right-hand side above. The diagonal elements can be obtained from the definition $M_{sk+j} \equiv \hat{T}_{sk+j} R_{sk+j} - R_{sk+j} \hat{T}_{sk+j}$, i.e.,

$$m_{1,1} = -\hat{\beta}_1 \rho_{1,2}, \qquad m_{sk+j+1,sk+j+1} = \hat{\beta}_{sk+j} \rho_{sk+j,sk+j+1}, \quad \text{and}$$
$$m_{i,i} = \hat{\beta}_{i-1} \rho_{i-1,i} - \hat{\beta}_i \rho_{i,i+1}, \quad \text{for} \quad i \in \{2, \ldots, sk+j\}.$$

Therefore, we can write

$$M_{sk+j} = \hat{T}_{sk+j} R_{sk+j} - R_{sk+j} \hat{T}_{sk+j} = \hat{\beta}_{sk+j+1} \hat{V}_{sk+j}^T \hat{v}_{sk+j+1} e_{sk+j+1}^T + H_{sk+j},$$

where $H_{sk+j}$ has elements satisfying

$$\begin{aligned}
\eta_{1,1} &= -\hat{\beta}_1 \rho_{1,2}, \\
\eta_{i,i} &= \hat{\beta}_{i-1} \rho_{i-1,i} - \hat{\beta}_i \rho_{i,i+1}, \quad \text{for} \quad i \in \{2, \ldots, sk+j+1\}, \\
\eta_{i-1,i} &= \hat{v}_{i-2}^T \delta\hat{v}_{i-1} - \delta\hat{v}_{i-2}^T \hat{v}_{i-1} + \hat{\beta}_{i-1}(\hat{v}_{i-2}^T \hat{v}_{i-2} - \hat{v}_{i-1}^T \hat{v}_{i-1}), \quad \text{and} \\
\eta_{\ell,i} &= \hat{v}_{\ell-1}^T \delta\hat{v}_{i-1} - \delta\hat{v}_{\ell-1}^T \hat{v}_{i-1}, \quad \text{for} \quad \ell \in \{1, \ldots, i-2\}.
\end{aligned} \qquad (4.35)$$

To simplify notation, we introduce the quantities

$$\bar{u}_{sk+j} = \max_{i \in \{0,\ldots,sk+j\}} \|\hat{u}_i\|_2, \quad \bar{\Gamma}_k = \max_{i \in \{0,\ldots,k\}} \Gamma_i, \quad \text{and} \quad \bar{\tau}_k = \max_{i \in \{0,\ldots,k\}} \tau_i.$$

Using this notation and (4.3), (4.23), (4.26), and (4.33), the quantities in (4.35) can be bounded by

$$|\eta_{1,1}| \leq 2\epsilon(n+11s+15)\bar{\Gamma}_k^2 \bar{u}_{sk+j}, \quad \text{and, for } i \in \{2, \ldots, sk+j+1\},$$

$$|\eta_{i,i}| \leq 4\epsilon(n+11s+15)\bar{\Gamma}_k^2 \bar{u}_{sk+j},$$

$$|\eta_{i-1,i}| \leq 2\epsilon\Big(\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k\big)\|A\|_2 + (n+18s+28)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2, \quad \text{and} \quad (4.36)$$

$$|\eta_{\ell,i}| \leq 2\epsilon\Big(\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k\big)\|A\|_2 + (10s+16)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2,$$

for $\ell \in \{1, \ldots, i-2\}$.

The above is a start toward proving (4.6). We return to this bound later, and now shift our focus towards proving a bound on $\|\hat{u}_{sk+j}\|_2$. To proceed, we must first find a bound for $|\hat{v}_{sk+j}^T \hat{v}_{sk+j-2}|$. From the definition of $M_{sk+j}$, we know the $(1,2)$ element of $M_{sk+j}$ is

$$\hat{\alpha}_0 \rho_{1,2} - \hat{\alpha}_1 \rho_{1,2} - \hat{\beta}_2 \rho_{1,3} = \eta_{1,2},$$

and for $i > 2$, the $(i-1, i)$ element is

$$\hat{\beta}_{i-2} \rho_{i-2,i} + (\hat{\alpha}_{i-2} - \hat{\alpha}_{i-1})\rho_{i-1,i} - \hat{\beta}_i \rho_{i-1,i+1} = \eta_{i-1,i}.$$

Then, defining

$$\xi_i \equiv (\hat{\alpha}_{i-2} - \hat{\alpha}_{i-1})\hat{\beta}_{i-1}\rho_{i-1,i} - \hat{\beta}_{i-1}\eta_{i-1,i}$$

for $i \in \{2, \ldots, sk+j\}$, we have

$$\hat{\beta}_{i-1}\hat{\beta}_i \rho_{i-1,i+1} = \hat{\beta}_{i-2}\hat{\beta}_{i-1}\rho_{i-2,i} + \xi_i = \xi_i + \xi_{i-1} + \ldots + \xi_2.$$

This, along with (4.19), (4.23), (4.26), and (4.36) gives

$$\hat{\beta}_{sk+j-1}\hat{\beta}_{sk+j}|\rho_{sk+j-1,sk+j+1}| = \hat{\beta}_{sk+j-1}\hat{\beta}_{sk+j}|\hat{v}_{sk+j-2}^T \hat{v}_{sk+j}|$$

$$\leq \sum_{i=2}^{sk+j} |\xi_i| \leq \sum_{i=2}^{sk+j} (|\hat{\alpha}_{i-2}| + |\hat{\alpha}_{i-1}|)|\hat{\beta}_{i-1}\rho_{i-1,i}| + |\hat{\beta}_{i-1}||\eta_{i-1,i}|$$

$$\leq 2\epsilon \sum_{i=2}^{sk+j} \Big(\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k\big)\|A\|_2 + (3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2 \bar{u}_{sk+j}$$

$$\leq 2\epsilon(sk+j-1)\Big(\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k\big)\|A\|_2 + (3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2 \bar{u}_{sk+j}.$$

$$(4.37)$$

Rearranging (4.30) gives

$$\hat{u}_{sk+j} - \delta u_{sk+j} = A\hat{v}_{sk+j} - \hat{\beta}_{sk+j}\hat{v}_{sk+j-1},$$

and multiplying each side by its own transpose (ignoring $\epsilon^2$ terms), we obtain

$$\hat{u}_{sk+j}^T \hat{u}_{sk+j} - 2\hat{u}_{sk+j}^T \delta u_{sk+j} = \|A\hat{v}_{sk+j}\|_2^2 + \hat{\beta}_{sk+j}^2 \|\hat{v}_{sk+j-1}\|_2^2 - 2\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T A\hat{v}_{sk+j-1}.$$

$$(4.38)$$

Rearranging (4.32) gives

$$A\hat{v}_{sk+j-1} = \hat{\beta}_{sk+j}\hat{v}_{sk+j} + \hat{\alpha}_{sk+j-1}\hat{v}_{sk+j-1} + \hat{\beta}_{sk+j-1}\hat{v}_{sk+j-2} - \delta\hat{v}_{sk+j-1},$$

and premultiplying this expression by $\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T$, we get

$$
\begin{aligned}
&\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T A\hat{v}_{sk+j-1}\\
&=\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T\big(\hat{\beta}_{sk+j}\hat{v}_{sk+j}+\hat{\alpha}_{sk+j-1}\hat{v}_{sk+j-1}+\hat{\beta}_{sk+j-1}\hat{v}_{sk+j-2}-\delta\hat{v}_{sk+j-1}\big)\\
&=\hat{\beta}_{sk+j}^2\|\hat{v}_{sk+j}\|_2^2+\hat{\alpha}_{sk+j-1}(\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T\hat{v}_{sk+j-1})+\hat{\beta}_{sk+j}\hat{\beta}_{sk+j-1}\hat{v}_{sk+j}^T\hat{v}_{sk+j-2}\\
&\quad-\hat{\beta}_{sk+j}\hat{v}_{sk+j}^T\delta\hat{v}_{sk+j-1}\\
&\equiv\hat{\beta}_{sk+j}^2+\delta\hat{\beta}_{sk+j},
\end{aligned}\tag{4.39}
$$

where, using bounds in (4.3), (4.19), (4.23), (4.26), (4.33), and (4.37),

$$
|\delta\hat{\beta}_{sk+j}|\le\epsilon\Big(\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k\big)\|A\|_2+(3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2\bar{u}_{sk+j}\tag{4.40}
$$
$$
+2\epsilon(sk+j-1)\Big(\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k\big)\|A\|_2+(3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2\bar{u}_{sk+j}.
$$

Adding $2\hat{u}_{sk+j}^T\delta u_{sk+j}$ to both sides of (4.38) and using (4.39), we obtain

$$
\begin{aligned}
\|\hat{u}_{sk+j}\|_2^2&=\|A\hat{v}_{sk+j}\|_2^2+\hat{\beta}_{sk+j}^2\big(\|\hat{v}_{sk+j-1}\|_2^2-2\big)-2\delta\hat{\beta}_{sk+j}+2\hat{u}_{sk+j}^T\delta u_{sk+j}\\
&\equiv\|A\hat{v}_{sk+j}\|_2^2+\hat{\beta}_{sk+j}^2\big(\|\hat{v}_{sk+j-1}\|_2^2-2\big)+\delta\tilde{\beta}_{sk+j},
\end{aligned}\tag{4.41}
$$

where $\delta\tilde{\beta}_{sk+j}=-2\delta\hat{\beta}_{sk+j}+2\hat{u}_{sk+j}^T\delta u_{sk+j}$, and, using the bounds in (4.31) and (4.40),

$$
\begin{aligned}
|\delta\tilde{\beta}_{sk+j}|&\le2|\delta\hat{\beta}_{sk+j}|+2\|\hat{u}_{sk+j}^T\|_2\|\delta u_{sk+j}\|_2\\
&\le4\epsilon(sk+j)\Big(\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k\big)\|A\|_2+(3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2\bar{u}_{sk+j}\\
&\quad+2\epsilon\Big(\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k\big)\|A\|_2+(4s+6)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2\bar{u}_{sk+j}\\
&\le4\epsilon(sk+j+1)\Big(\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k\big)\|A\|_2+(3n+40s+58)\bar{u}_{sk+j}\Big)\bar{\Gamma}_k^2\bar{u}_{sk+j}.
\end{aligned}\tag{4.42}
$$

Now, using (4.41), and since $\hat{\beta}_{sk+j}^2\ge0$, we can write

$$
\|\hat{u}_{sk+j}\|_2^2\le\|\hat{u}_{sk+j}\|_2^2+\hat{\beta}_{sk+j}^2\le\|A\|_2^2\|\hat{v}_{sk+j}\|_2^2+\hat{\beta}_{sk+j}^2\big(\|\hat{v}_{sk+j-1}\|_2^2-1\big)+|\delta\tilde{\beta}_{sk+j}|.\tag{4.43}
$$

Let $\mu\equiv\max\big\{\bar{u}_{sk+j},\|A\|_2\big\}$. Then (4.43) along with bounds in (4.3), (4.23), and (4.42) gives

$$
\|\hat{u}_{sk+j}\|_2^2\le\|A\|_2^2+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\mu^2.\tag{4.44}
$$

We consider the two possible cases for $\mu$. First, if $\mu=\|A\|_2$, then

$$
\begin{aligned}
\|\hat{u}_{sk+j}\|_2^2&\le\|A\|_2^2+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\|A\|_2^2\\
&\le\|A\|_2^2\Big(1+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\Big).
\end{aligned}
$$

Otherwise, we have the case $\mu=\bar{u}_{sk+j}$. Since the bound in (4.44) holds for all $\|\hat{u}_{sk+j}\|_2^2$, it also holds for $\bar{u}_{sk+j}^2=\mu^2$, and thus, ignoring terms of order $\epsilon^2$,

$$
\begin{aligned}
\mu^2&\le\|A\|_2^2+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\mu^2\\
&\le\|A\|_2^2+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\|A\|_2^2\\
&\le\|A\|_2^2\Big(1+4\epsilon(sk+j+2)\big((n+2s+5)\sigma+(4s+9)\bar{\tau}_k+(3n+40s+58)\big)\bar{\Gamma}_k^2\Big),
\end{aligned}
$$

and, plugging this in to (4.44), we get

$$\|\hat{u}_{sk+j}\|_2^2 \le \|A\|_2^2 \Big(1 + 4\epsilon(sk+j+2)\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (3n+40s+58)\big)\bar{\Gamma}_k^2\Big). \quad (4.45)$$

In either case we obtain the same bound on $\|\hat{u}_{sk+j}\|_2^2$, so (4.45) holds.

Taking the square root of (4.45), we have

$$\|\hat{u}_{sk+j}\|_2 \le \|A\|_2 \Big(1 + 2\epsilon(sk+j+2)\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (3n+40s+58)\big)\bar{\Gamma}_k^2\Big), \quad (4.46)$$

and substituting (4.46) into (4.26), (4.33), and (4.36) proves the bounds (4.2), (4.1), and (4.6) in Theorem 4.2, respectively, assuming that

$$2\epsilon(sk+j+2)\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (3n+40s+58)\big)\bar{\Gamma}_k^2 \ll 1.$$

The only remaining inequality to prove is (4.4). To do this, we first multiply both sides of (4.24) by their own transposes to obtain

$$\hat{\beta}_{sk+j+1}^2 \|\hat{v}_{sk+j+1}\|_2^2 = \|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + 2\delta w_{sk+j}^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}'.$$

Adding $\hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2$ to both sides,

$$\hat{\beta}_{sk+j+1}^2 \|\hat{v}_{sk+j+1}\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 = \|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 \\ + 2\delta w_{sk+j}^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}'.$$

Substituting in (4.41) on the left hand side,

$$\hat{\beta}_{sk+j+1}^2 \|\hat{v}_{sk+j+1}\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|A\hat{v}_{sk+j}\|_2^2 - \hat{\beta}_{sk+j}^2(\|\hat{v}_{sk+j-1}\|_2^2 - 2) - \delta\tilde{\beta}_{sk+j} = \\ \|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 + 2\delta w_{sk+j}^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}',$$

and then subtracting $\hat{\beta}_{sk+j+1}^2$ from both sides gives

$$\hat{\beta}_{sk+j+1}^2(\|\hat{v}_{sk+j+1}\|_2^2 - 1) + \hat{\alpha}_{sk+j}^2 - \|A\hat{v}_{sk+j}\|_2^2 - \hat{\beta}_{sk+j}^2(\|\hat{v}_{sk+j-1}\|_2^2 - 2) - \delta\tilde{\beta}_{sk+j} = \\ \|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 + 2\delta w_{sk+j}^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}' - \hat{\beta}_{sk+j+1}^2.$$

This can be rearranged to give

$$\hat{\beta}_{sk+j+1}^2 + \hat{\alpha}_{sk+j}^2 + \hat{\beta}_{sk+j}^2 - \|A\hat{v}_{sk+j}\|_2^2 = \|\hat{\mathcal{Y}}_k \hat{w}_{k,j}'\|_2^2 + \hat{\alpha}_{sk+j}^2 - \|\hat{u}_{sk+j}\|_2^2 \\ + 2\delta w_{sk+j}^T \hat{\mathcal{Y}}_k \hat{w}_{k,j}' + \hat{\beta}_{sk+j}^2(\|\hat{v}_{sk+j-1}\|_2^2 - 1) \\ - \hat{\beta}_{sk+j+1}^2(\|\hat{v}_{sk+j+1}\|_2^2 - 1) + \delta\tilde{\beta}_{sk+j},$$

and finally, using (4.3), (4.21), (4.23), (4.25), and (4.42) gives the bound

$$\left| \hat{\beta}_{sk+j+1}^2 + \hat{\alpha}_{sk+j}^2 + \hat{\beta}_{sk+j}^2 - \|A\hat{v}_{sk+j}\|_2^2 \right| \le \\ 4\epsilon(sk+j+2)\big((n+2s+5)\sigma + (4s+9)\bar{\tau}_k + (3n+40s+58)\big)\bar{\Gamma}_k^2 \|A\|_2^2.$$

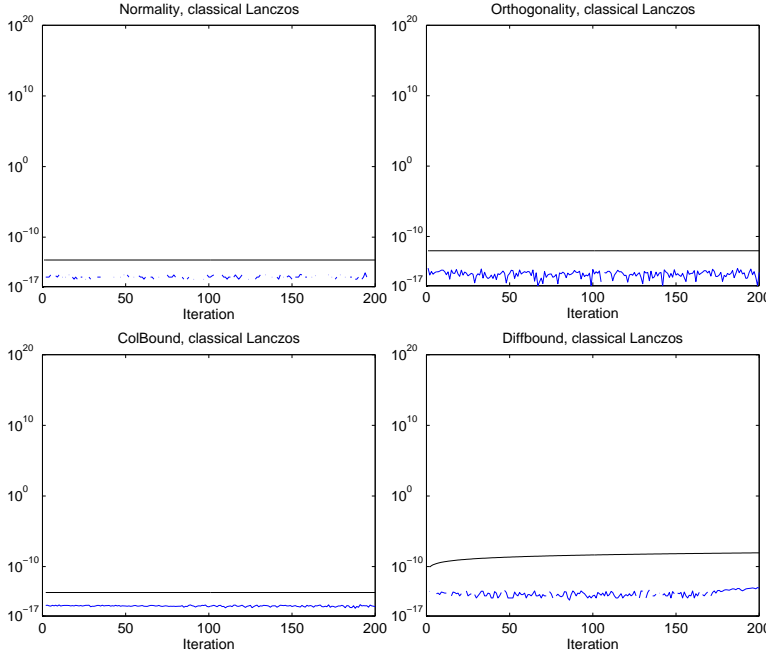This proves (4.4) and thus completes the proof of Theorem 4.2.

FIG. 5.1. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for classical Lanczos on 2D Poisson with $n = 256$. Upper bounds are taken from [27].*

**5. Numerical Examples.** We give a brief example to illustrate the bounds in (4.1), (4.2), (4.3), and (4.4). We run $s$-step Lanczos (Algorithm 2) in double precision with $s = 8$ on the same model problem used in Section 3: a 2D Poisson matrix with $n = 256$, $\|A\|_2 = 7.93$, using a random starting vector. For comparison, Figure 5.1 shows the results for classical Lanczos using the bounds derived by Paige [27]. In the top left, the blue curve gives the measured value of normality, $|\hat{v}_{i+1}^T \hat{v}_{i+1} - 1|$, and the black curve plots the upper bound, $(n+4)\epsilon$. In the top right, the blue curve gives the measured value of orthogonality, $|\hat{\beta}_{i+1} \hat{v}_i^T \hat{v}_{i+1}|$, and the black curve plots the upper bound, $2(n+4)\|A\|_2\epsilon$. In the bottom left, the blue curve gives the measured value of the bound (4.1) for $\|\delta\hat{v}_i\|_2$, and the black curve plots the upper bound, $\epsilon(7+5\|\,|A|\,\|_2)$. In the bottom right, the blue curve gives the measured value of the bound (4.4), and the black curve plots the upper bound, $4i\epsilon(3(n+4)\|A\|_2 + (7+5\|\,|A|\,\|_2))\|A\|_2$.

The results for $s$-step Lanczos are shown in Figures 5.2−5.4. The same tests were run for three different basis choices: monomial (Figure 5.2), Newton (Figure 5.3), and Chebyshev (Figure 5.4) (see, e.g., [31]). For each of the four plots in each Figure, the blue curves give the measured values of the quantities on the left hand sides of (clockwise from the upper left) (4.3), (4.2), (4.1), and (4.4). The cyan curves give the maximum of the measured values so far. The red curves give the value of $\bar{\Gamma}_k^2$ as defined in Theorem 4.2, and the blacks curves give the upper bounds on the right hand sides of (4.3), (4.2), (4.1), and (4.4).

We see from Figures 5.2−5.4 that the upper bounds given in Theorem 4.2 are valid. In particular, we can also see that the shape of the curve $\bar{\Gamma}_k^2$ gives a good indication of the shape of the curves for $\max_{i \le sk+j} |\hat{v}_{i+1}^T \hat{v}_{i+1} - 1|$ and $\max_{i \le sk+j} |\hat{\beta}_{i+1} \hat{v}_i^T \hat{v}_{i+1}|$. However, from Figure 5.2 for the monomial basis, we see that if the basis has a high
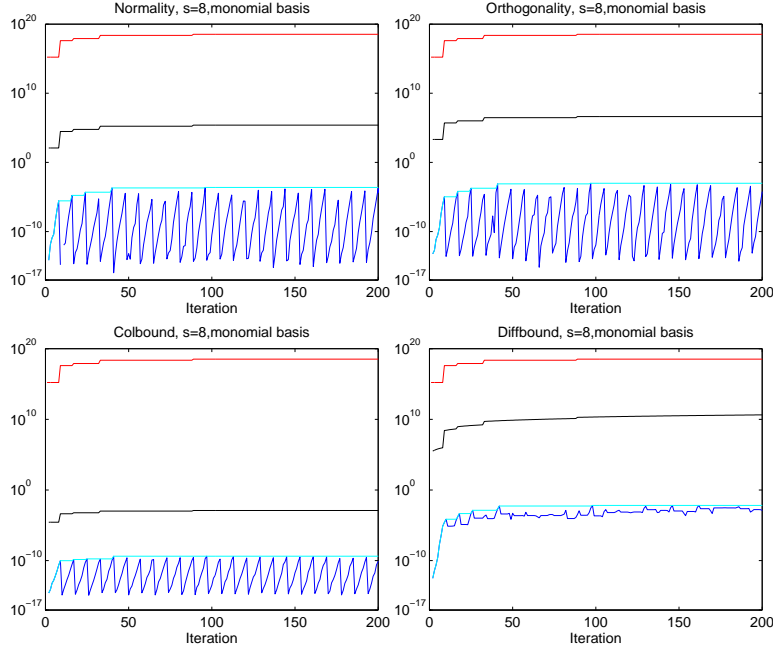
Fig. 5.2. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for s-step Lanczos on 2D Poisson with $n = 256$ and $s = 8$ for monomial basis. Bounds obtained using $\bar{\Gamma}_k$ as defined in Theorem 4.2.*
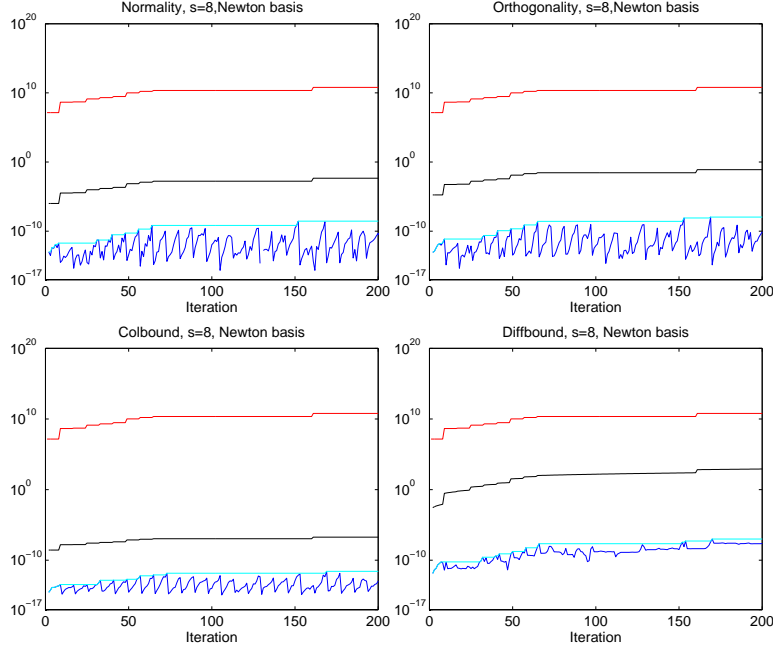


Fig. 5.3. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for s-step Lanczos on 2D Poisson with $n = 256$ and $s = 8$ for Newton basis. Bounds obtained using $\bar{\Gamma}_k$ as defined in Theorem 4.2.*
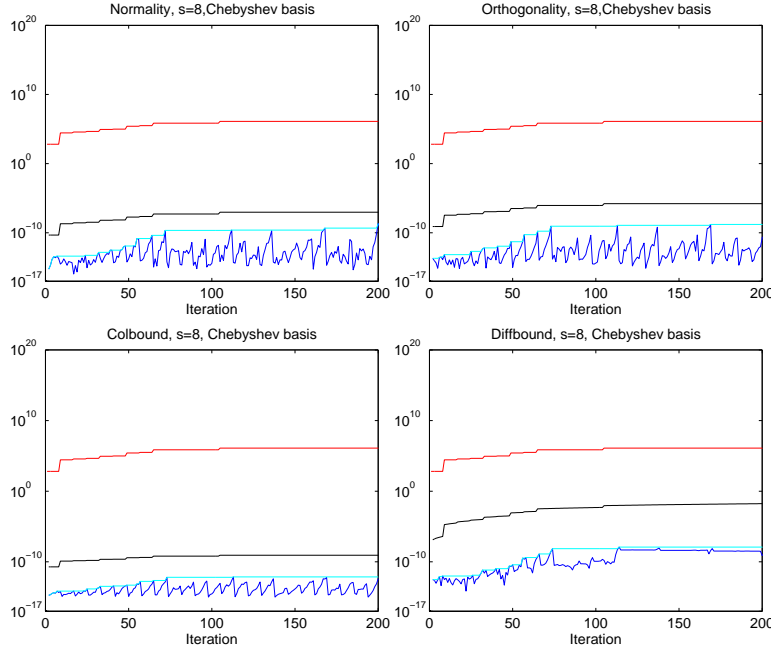
FIG. 5.4. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for s-step Lanczos on 2D Poisson with $n = 256$ and $s = 8$ for Chebyshev basis. Bounds obtained using $\bar{\Gamma}_k$ as defined in Theorem 4.2.*

condition number, as does the monomial basis here, the upper bound can be a very large overestimate quantitatively, leading to bounds that are not useful.

There is an easy way to improve the bounds by using a different definition of $\bar{\Gamma}_k$ to upper bound quantities in the proof of Theorem 4.2. Note that all quantities which we have bounded by $\bar{\Gamma}_k$ in Section 4 are of the form $\| |\hat{\mathcal{Y}}_k| |x| \|_2 / \|\hat{\mathcal{Y}}_k x\|_2$. While the use of $\bar{\Gamma}_k$ as defined in Theorem 4.2 shows how the bounds depend on the conditioning of the computed Krylov bases, we can obtain tighter and more descriptive bounds for (4.3) and (4.2) by instead using the definition

$$\bar{\Gamma}_{k,j} \equiv \max_{x \in \{\hat{w}'_{k,j}, \hat{u}'_{k,j}, \hat{v}'_{k,j}, \hat{v}'_{k,j-1}\}} \frac{\| |\hat{\mathcal{Y}}_k| |x| \|_2}{\|\hat{\mathcal{Y}}_k x\|_2}. \tag{5.1}$$

For the bound in (4.1), we can use the definition

$$\bar{\Gamma}_{k,j} \equiv \max \left\{ \frac{\| |\hat{\mathcal{Y}}_k| |\mathcal{B}_k| |\hat{v}'_{k,j}| \|_2}{\| |\mathcal{B}_k| \|_2 \|\hat{\mathcal{Y}}_k \hat{v}'_{k,j}\|_2}, \max_{x \in \{\hat{w}'_{k,j}, \hat{u}'_{k,j}, \hat{v}'_{k,j}, \hat{v}'_{k,j-1}\}} \frac{\| |\hat{\mathcal{Y}}_k| |x| \|_2}{\|\hat{\mathcal{Y}}_k x\|_2} \right\}, \tag{5.2}$$

and for the bound in (4.4), we can use the definition

$$\bar{\Gamma}_{k,j} \equiv \max \left\{ \bar{\Gamma}_{k,j-1}, \frac{\| |\hat{\mathcal{Y}}_k| |\mathcal{B}_k| |\hat{v}'_{k,j}| \|_2}{\| |\mathcal{B}_k| \|_2 \|\hat{\mathcal{Y}}_k \hat{v}'_{k,j}\|_2}, \max_{x \in \{\hat{w}'_{k,j}, \hat{u}'_{k,j}, \hat{v}'_{k,j}, \hat{v}'_{k,j+1}\}} \frac{\| |\hat{\mathcal{Y}}_k| |x| \|_2}{\|\hat{\mathcal{Y}}_k x\|_2} \right\}. \tag{5.3}$$

The value in (5.3) is monotonically increasing since the bound in (4.37) is a sum of bounds from previous iterations.
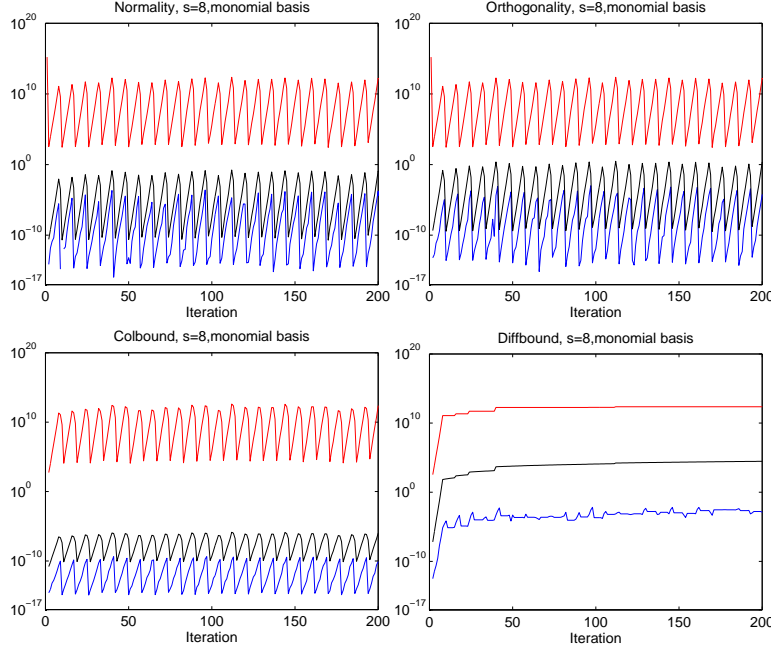
FIG. 5.5. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for monomial basis. Bounds obtained using $\bar{\Gamma}_k$ as defined in* (5.1) *for top plots,* (5.2) *for bottom left plot, and* (5.3) *for bottom right plot.*

In Figures 5.5−5.7 we plot bounds for the same problem, bases, and $s$-values as Figures 5.2−5.4, but using the new definitions of $\bar{\Gamma}_{k,j}$. Comparing Figures 5.5−5.7 to Figures 5.2−5.4, we see that these bounds are better both quantitatively, in that they are tighter, and qualitatively, in that they better replicate the shape of the curves for the measured normality and orthogonality values. The exception is for the plots of bounds in (4.4) (bottom right plots), for which there is not much difference qualitatively. It is also clear that the new definitions of $\bar{\Gamma}_k$ correlate well with the size of the measured values (i.e., the shape of the blue curve closely follows the shape of the red curve). Note that, unlike the definition of $\bar{\Gamma}_k$ in Theorem 4.2, using the definitions in (5.1)−(5.3) do not require the assumption of linear independence of the basis vectors.

Although these new bounds can not be computed a priori, the right hand sides of (5.1), (5.2), and (5.3) can be computed within each inner loop iteration for the cost of one extra reduction per outer loop. This extra cost comes from the need to compute $|\hat{\mathcal{Y}}_k|^T|\hat{\mathcal{Y}}_k|$, although this could potentially be performed simultaneously with the computation of $\hat{G}_k$ (line 4 in Algorithm 2). This means that meaningful bounds could be cheaply estimated during the iterations. Designing a scheme to improve numerical properties using this information remains future work.

**6. Future work.** In this paper, we have presented a complete rounding error analysis of the $s$-step Lanczos method. The derived bounds are analogous to those of Paige for classical Lanczos [27], but also depend on a amplification factor $\bar{\Gamma}_k^2$, which depends on the condition number of the Krylov bases computed every in each outer loop. Our analysis confirms the empirical observation that the conditioning of the
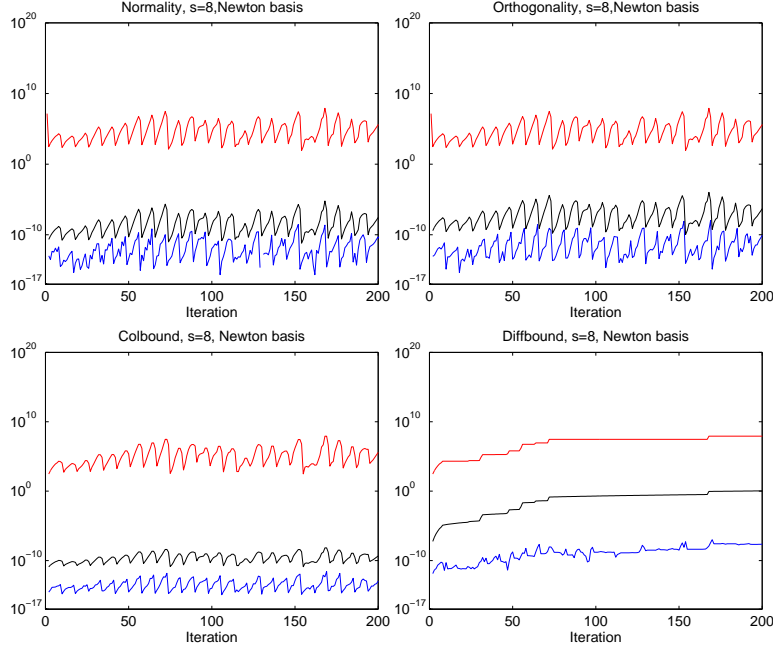
Fig. 5.6. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for Newton basis. Bounds obtained using* $\bar{\Gamma}_k$ *as defined in* (5.1) *for top plots,* (5.2) *for bottom left plot, and* (5.3) *for bottom right plot.*
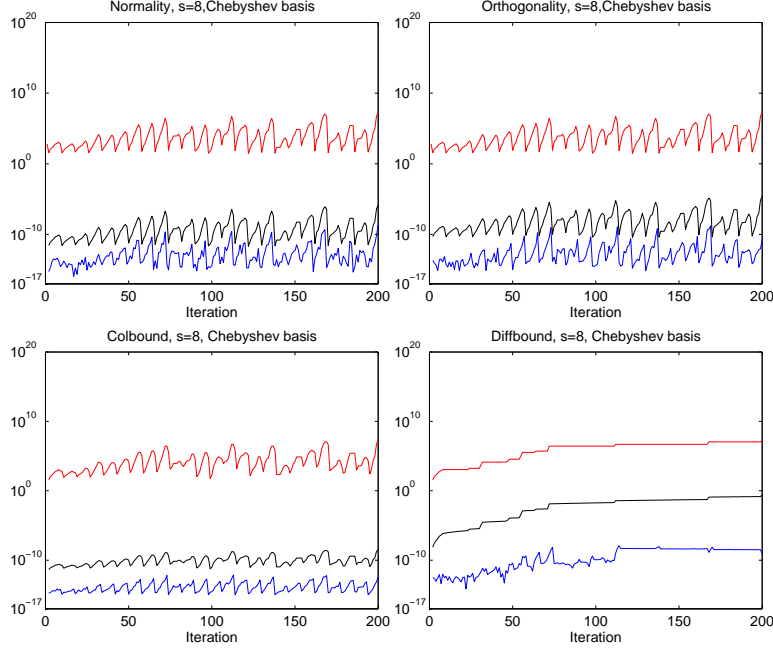


Fig. 5.7. *Normality (top left), orthogonality (top right), recurrence column error (bottom left), and difference in recurrence column size (bottom right) for Chebyshev basis. Bounds obtained using* $\bar{\Gamma}_k$ *as defined in* (5.1) *for top plots,* (5.2) *for bottom left plot, and* (5.3) *for bottom right plot.*

Krylov bases plays a large role in determining finite precision behavior.

The next step is to extend the analogous subsequent analyses of Paige, in which he proves properties about Ritz vectors and Ritz values, relates the convergence of a Ritz pair to loss of orthogonality, and, more recently, proves a type of augmented backward stability for the classical Lanczos method [28, 29].

Another area of interest is the development of practical techniques for improving $s$-step Lanczos based on our results. This could include strategies for reorthogonalizing the Lanczos vectors, (re)orthogonalizing the generated Krylov basis vectors, or controlling the basis conditioning in a number of ways. The bounds could also be used for guiding the use of extended precision in $s$-step Krylov methods; for example, if we want the bounds in Theorem 4.2 for the $s$-step method with precision $\tilde{\epsilon}$ to be similar to those for the classical method with precision $\epsilon$, one must use precision $\tilde{\epsilon} \approx \epsilon/\bar{\Gamma}_k^2$.

In this analysis, our upper bounds are likely large overestimates. This is in part due to our replacing $\Gamma_k$ with $\Gamma_k^2$ in order to simplify many of the bounds. If the analysis in this paper is performed instead keeping both $\Gamma_k$ and $\Gamma_k^2$ terms, it can be shown that increasing the precision in a few computations (involving the construction and application of the Gram matrix $\hat{G}_k$) can improve the error bounds in Theorem 4.2 by a factor of $\bar{\Gamma}_k$. This motivates the development of mixed precision $s$-step Lanczos methods, which could potentially trade bandwidth (in extra bits of precision) for fewer total iterations. As demonstrated in Section 5, it is also possible to use a tighter, iteratively updated bound for $\bar{\Gamma}_k$ which results in tighter and more descriptive bounds for the quantities in Theorem 4.2.

## REFERENCES

[1] Z. BAI, D. HU, AND L. REICHEL, *A Newton basis GMRES implementation*, IMA J. Numer. Anal., 14 (1994), pp. 563–581.

[2] G. BALLARD, E. CARSON, J. DEMMEL, M. HOEMMEN, N. KNIGHT, AND O. SCHWARTZ, *Communication lower bounds and optimal algorithms for numerical linear algebra*, Acta Numer. (in press), (2014).

[3] E. CARSON AND J. DEMMEL, *A residual replacement strategy for improving the maximum attainable accuracy of s-step Krylov subspace methods*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 22–43.

[4] E. CARSON, N. KNIGHT, AND J. DEMMEL, *Avoiding communication in nonsymmetric Lanczos-based Krylov subspace methods*, SIAM J. Sci. Comp., 35 (2013).

[5] A. CHRONOPOULOS AND C. GEAR, *On the efficient implementation of preconditioned s-step conjugate gradient methods on multiprocessors with memory hierarchy*, Parallel Comput., 11 (1989), pp. 37–53.

[6] ———, *s-step iterative methods for symmetric linear systems*, J. Comput. Appl. Math, 25 (1989), pp. 153–168.

[7] A. CHRONOPOULOS AND C. SWANSON, *Parallel iterative s-step methods for unsymmetric linear systems*, Parallel Comput., 22 (1996), pp. 623–641.

[8] E. DE STURLER, *A performance model for Krylov subspace methods on mesh-based parallel computers*, Parallel Comput., 22 (1996), pp. 57–74.

[9] J. DEMMEL, M. HOEMMEN, M. MOHIYUDDIN, AND K. YELICK, *Avoiding communication in computing Krylov subspaces*, Tech. Report UCB/EECS-2007-123, EECS Dept., U.C. Berkeley, Oct 2007.

[10] D. GANNON AND J. VAN ROSENDALE, *On the impact of communication complexity on the design of parallel numerical algorithms*, Trans. Comput., 100 (1984), pp. 1180–1194.

[11] G. GOLUB AND C. VAN LOAN, *Matrix computations*, JHU Press, Baltimore, MD, 3 ed., 1996.

[12] A. GREENBAUM AND Z. STRAKOŠ, *Predicting the behavior of finite precision Lanczos and conjugate gradient computations*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 121–137.

[13] M. GUSTAFSSON, J. DEMMEL, AND S. HOLMGREN, *Numerical evaluation of the communication-avoiding Lanczos algorithm*, Tech. Report ISSN 1404-3203/2012-001, Department of Information Technology, Uppsala University, Feb. 2012.

[14] M. GUTKNECHT, *Lanczos-type solvers for nonsymmetric linear systems of equations*, Acta Numer., 6 (1997), pp. 271–398.

[15] M. GUTKNECHT AND Z. STRAKOŠ, *Accuracy of two three-term and three two-term recurrences for Krylov space solvers*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 213–229.

[16] A. HINDMARSH AND H. WALKER, *Note on a Householder implementation of the GMRES method*, Tech. Report UCID-20899, Lawrence Livermore National Lab., CA., 1986.

[17] M. HOEMMEN, *Communication-avoiding Krylov subspace methods*, PhD thesis, EECS Dept., U.C. Berkeley, 2010.

[18] W. JOUBERT AND G. CAREY, *Parallelizable restarted iterative methods for nonsymmetric linear systems. Part I: theory*, Int. J. Comput. Math., 44 (1992), pp. 243–267.

[19] S. KIM AND A. CHRONOPOULOS, *A class of Lanczos-like algorithms implemented on parallel computers*, Parallel Comput., 17 (1991), pp. 763–778.

[20] ———, *An efficient nonsymmetric Lanczos method on parallel vector computers*, J. Comput. Appl. Math., 42 (1992), pp. 357–374.

[21] C. LANCZOS, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators1*, J. Res. Natn. Bur. Stand., 45 (1950), pp. 255–282.

[22] G. MEURANT, *The Lanczos and conjugate gradient algorithms: from theory to finite precision computations*, SIAM, 2006.

[23] G. MEURANT AND Z. STRAKOŠ, *The Lanczos and conjugate gradient algorithms in finite precision arithmetic*, Acta Numer., 15 (2006), pp. 471–542.

[24] M. MOHIYUDDIN, M. HOEMMEN, J. DEMMEL, AND K. YELICK, *Minimizing communication in sparse matrix solvers*, in Proc. ACM/IEEE Conference on Supercomputing, 2009.

[25] C. PAIGE, *The computation of eigenvalues and eigenvectors of very large sparse matrices*, PhD thesis, London University, London, UK, 1971.

[26] ———, *Computational variants of the Lanczos method for the eigenproblem*, IMA J. Appl. Math., 10 (1972), pp. 373–381.

[27] ———, *Error analysis of the Lanczos algorithm for tridiagonalizing a symmetric matrix*, IMA J. Appl. Math., 18 (1976), pp. 341–349.

[28] ———, *Accuracy and effectiveness of the Lanczos algorithm for the symmetric eigenproblem*, Linear Algebra Appl., 34 (1980), pp. 235–258.

[29] ———, *An augmented stability result for the Lanczos hermitian matrix tridiagonalization process*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2347–2359.

[30] B. PARLETT AND D. SCOTT, *The Lanczos algorithm with selective orthogonalization*, Math. Comput., 33 (1979), pp. 217–238.

[31] B. PHILIPPE AND L. REICHEL, *On the generation of Krylov subspace bases*, Appl. Numer. Math, 62 (2012), pp. 1171–1186.

[32] H. SIMON, *The Lanczos algorithm with partial reorthogonalization*, Math. Comput., 42 (1984), pp. 115–142.

[33] S. TOLEDO, *Quantitative performance modeling of scientific computations and creating locality in numerical algorithms*, PhD thesis, MIT, 1995.

[34] H. VAN DER VORST AND Q. YE, *Residual replacement strategies for Krylov subspace iterative methods for the convergence of true residuals*, SIAM J. Sci. Comput., 22 (1999), pp. 835–852.

[35] J. VAN ROSENDALE, *Minimizing inner product data dependencies in conjugate gradient iteration*, Tech. Report 172178, ICASE-NASA, 1983.

[36] H. WALKER, *Implementation of the GMRES method using Householder transformations*, SIAM J. Sci. Stat. Comput., 9 (1988), pp. 152–163.

[37] W. WÜLLING, *On stabilization and convergence of clustered Ritz values in the Lanczos method*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 891–908.

[38] J. ZEMKE, *Krylov subspace methods in finite precision: a unified approach*, PhD thesis, Technische Universität Hamburg-Harburg, 2003.